

From Explainable AI to Human-Centered System Reliability: Quantifying and Visualizing Calibrated Trust in Mission-Critical XR

Matthew Wilchek^{1,*}, Kurt Luther² and Feras A. Batarseh³

¹U.S. Army, Combat Capabilities Development Command, C5ISR Center, Fort Belvoir, VA, USA

²Department of Computer Science and Center for HCI, Virginia Tech, Alexandria, VA, USA

³Department of Biological Systems Engineering, Virginia Tech, Arlington, VA, USA

Abstract

Extended Reality (XR) systems increasingly integrate Artificial Intelligence (AI) to support professionals in high-risk settings such as search and rescue, law enforcement, and military operations. Yet common approaches to trust and explainability often rely on qualitative assessments or offline explanations that are poorly suited to embodied, time-critical work. This paper outlines an emerging perspective on how mission-critical XR systems may better support calibrated trust, resilient oversight, and situated explainability through human-centered system reliability (HCSR), a quantitative, user-parameterized estimate of reliability that evolves through accumulated evidence. Drawing on prior work in distributed XR and human-AI teaming, we describe three connected shifts: from generic trust to calibrated trust through updated reliability estimates; from seamlessness to resilience through human-centered oversight; and from transparency to situated explainability through lightweight spatial cues embedded in the XR interface. We conclude with implications and open challenges for reliability-aware XR systems.

Keywords

Embodied Explainable AI, Trust Calibration, Mission-Critical Systems, Mixed Reality, Human-in-the-Loop AI, Shared Perception

1. Introduction

Extended Reality (XR) systems increasingly integrate Artificial Intelligence (AI) to support professionals working in high-risk settings such as search and rescue (SAR), law enforcement, and military operations. These contexts exemplify *extreme sensemaking*, requiring high-stakes decision-making under uncertainty, time pressure, and incomplete information in sparse environments [1]. Decades of military XR research have shown that head-up, spatially registered overlays can improve situational awareness (SA) by embedding relevant information directly into the user’s view of the physical environment [2]. However, this same body of work also shows that simply adding information can overwhelm users and degrade SA, which motivates role-aware filtering and careful cue presentation [3].

In parallel, trust and explainability are widely viewed as prerequisites for effective AI-assisted decision-making [4]. Early work conceptualized trust through qualitative dimensions such as perceived reliability, technical competence, and understandability [5, 6], while more recent surveys [7] emphasize the importance of explainable AI (XAI) in fostering user confidence, particularly in computer vision systems [8]. Yet these approaches are largely designed for offline inspection and reflection and are not well suited to the embodied, time-critical contexts XR users face in the field [9, 10]. Trust is a key determinant of team efficiency in human-AI collaboration, but its role is interpreted inconsistently across “trustworthy AI” literature [11, 12, 13]. In distributed teams, appropriately calibrated trust allows members to act on AI insights while reducing coordination delays and decision friction [14, 15].

CHI’26: XR4CE workshop, April 14, 2026, Barcelona, Spain

*Corresponding author

✉ matthew.r.wilchek.civ@army.mil (M. Wilchek); kluther@vt.edu (K. Luther); batarseh@vt.edu (F. A. Batarseh)

🆔 0000-0001-8664-1586 (M. Wilchek); 0000-0003-1809-6269 (K. Luther); 0000-0002-6062-2747 (F. A. Batarseh)



© 2026 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

Following Gunning [16], we treat trustworthy AI not as algorithms that report confidence, but as systems that communicate reliability, uncertainty, and provenance in ways that support human judgment.

The PerceptiSync framework builds on this foundation by providing a human-centered implementation of Dirichlet-Categorical (DC) trust modeling for distributed AI systems, incorporating user-configurable features and real-time human feedback to dynamically adjust trust assessments [17]. Empirical evaluations show that integrating Human-in-the-Loop (HITL) and Crowd-in-the-Loop (CITL) mechanisms can improve trust assessment and system reliability in dynamic environments [18, 19]. Despite these advances, more research is needed to understand the design and impacts of human-centered system reliability (HCSR) estimates implemented for XR systems. Mixed-reality systems for challenging environments often focus on what information should be shown and where, while trust frameworks focus on how reliable information sources are, with less attention to how those assessments are conveyed to embodied users. For this workshop, we organize that perspective around three connected shifts, as depicted in Figure 1. Shift 1 positions HCSR as a mechanism for dynamically calibrating trust through accumulated evidence. Shift 2 explains how these evolving reliability estimates can support resilient, human-centered oversight when uncertainty, disagreement, or potential breakdowns occur. Shift 3 positions situated explainability in XR as the interface layer that makes those reliability changes and intervention opportunities visible in real time. Together, these shifts outline an emerging view of how XR systems for mission-critical work can better support calibrated trust, resilient oversight, and situated explainability.

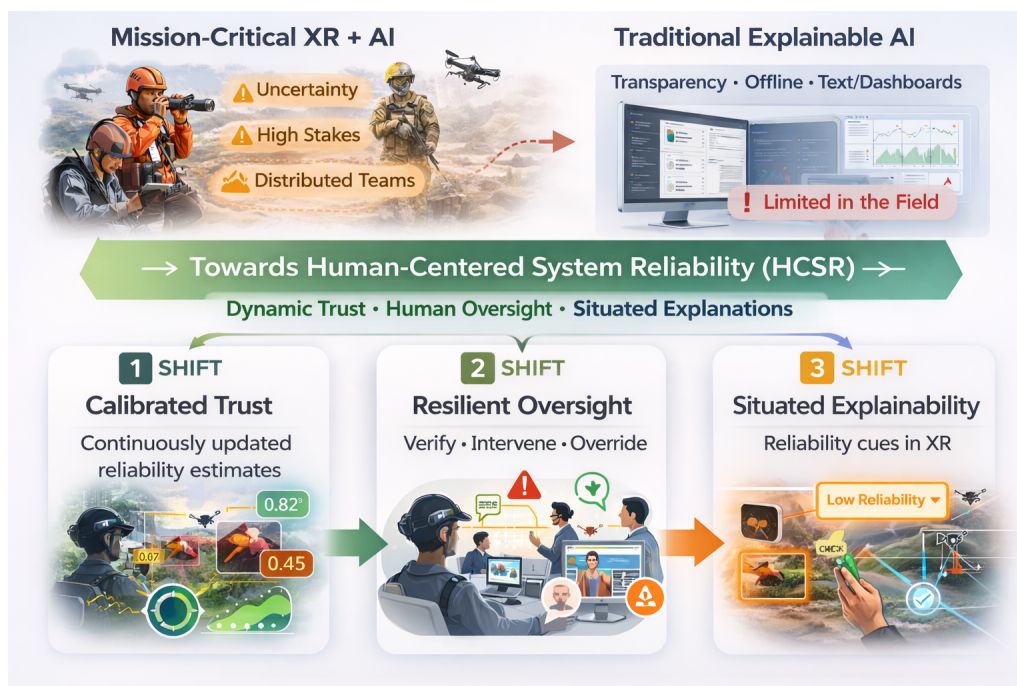


Figure 1: Moving from traditional explainable AI toward human-centered system reliability in mission-critical XR through three connected shifts: calibrated trust, resilient oversight, and situated explainability.

2. Shift 1: Calibrated Trust via Human-Centered Reliability

A common approach to studying trust in AI-assisted systems relies on post-hoc, subjective methods such as Likert-style questionnaires or self-reported scales [20]. While these measures provide insight into users' subjective perceptions, they offer limited support for real-time decision-making in dynamic, high-risk environments. In challenging environments, professionals must continually decide whether to rely on specific AI outputs, teammates, or sensing modalities as situations evolve. This motivates a move away from static trust scores and toward reliability as a measurable, continuously updated signal

that reflects how information sources perform over time.

Prior work has begun to formalize trust as a probabilistic, evidence-driven process. In particular, Guo et al. introduced a DC trust model in which trust between a source i and receiver j is represented as an opinion parameterized by accumulated positive, negative, and uncertain evidence [21]. The resulting reliability estimate is computed as:

$$\omega_{ij} = \frac{\alpha_{ij}}{\alpha_{ij} + \beta_{ij} + \gamma_{ij}} \quad (1)$$

where α_{ij} , β_{ij} , and γ_{ij} denote the amount of positive, negative, and uncertain evidence, respectively. As new observations are exchanged, these parameters are continuously updated, yielding a time-evolving estimate of source reliability. Building on this foundation, Cirne et al. show that such evolving trajectories exhibit recognizable patterns of growth, decay, or oscillation as evidence accumulates or conflicts over time [22]. Reliability estimates differ from both model confidence and explanation. Model confidence reflects an algorithm’s internal uncertainty about a specific prediction, while explanations aim to clarify why a prediction was produced. Reliability, by contrast, captures how dependable a source has been across interactions and time, independent of any single output. This distinction matters in distributed XR systems, where users must reason across heterogeneous AI agents, sensing modalities, and teammates whose performance may vary across contexts.

We extend this probabilistic reliability model into *human-centered system reliability* (HCSR), a user-parameterized estimate of an information source’s reliability that evolves through accumulated evidence exchange *and* individual trust preferences. The PerceptiSync framework implements this idea by incorporating user-configurable features that govern how conservatively or quickly reliability evolves, and by integrating real-time human feedback into trust updates [17]. Treating HCSR as a foundational mechanism enables calibrated trust by allowing XR systems to convey both what an AI detected and how much a user should rely on that source. This, in turn, supports resilient oversight and situated explainability.

3. Shift 2: Resilience through Human-Centered Oversight

Prior research on AI-enabled wearable and XR systems identifies trust, adaptability, cognitive burden, and error tolerance as central challenges for human-in-the-loop operation, noting that small errors or failures can present significant safety and performance risks in real-world deployments [19]. These findings suggest that mission-critical XR systems need to support human oversight when uncertainty arises [23]. In extreme sensemaking, resilience can be improved through intuitive user interfaces and HITL mechanisms that enable humans to perceive uncertainty, add contextual knowledge, and intervene when automated assessments diverge from reality. Within this framing, changing reliability estimates should not remain passive analytics in the background. They should actively inform when oversight is needed.

Our prior work on PerceptiSync builds on HCSR by providing explicit mechanisms for human-centered oversight grounded in CITL and HITL principles. Users can define a trust persona composed of multiple configuration features that determine when humans are consulted, when intervention is required, and how user judgment shapes reliability estimates.¹ These features include *trust level*, which specifies how conservative or permissive a system should be when incorporating new information; *trust history*, which determines whether accumulated past evidence must be present before reliability is adjusted; and *trustworthy assessments*, which specify how many consecutive reliable observations are required before trust is reinforced or degraded. Together, these parameters allow reliability estimates to reflect both accumulated evidence and individual user preferences.

These features are not fixed settings. Instead, they determine how incoming evidence changes HCSR over time and when the system should transition from autonomous reliance to human-centered oversight. After a new observation, detection, or message is received, PerceptiSync records the evidence as positive, negative, or uncertain, depending on whether the information is corroborated, contradicted,

¹The GitHub source for the PerceptiSync algorithm can be found [here](#).

or remains ambiguous. The user's trust persona then determines how strongly that evidence affects the current reliability estimate by shaping how conservative the update should be, how much prior history is required, and how many trusted events are needed before a source is reinforced or degraded.

In practice, dynamic adjustment can occur in two complementary ways. First, users may configure their trust persona based on mission context. For example, they may select a more cautious trust level during low-visibility SAR operations or when integrating a newly deployed sensing modality. Second, the system can adapt implicitly as evidence accumulates during interaction. Repeated agreement between distributed agents can strengthen reliability and reduce the need for intervention, while conflicting observations, uncertain detections, or missing corroboration can slow trust growth, lower reliability, or trigger oversight prompts. In this way, calibrated trust is adjusted through a continuous loop of evidence updates, user preferences, and possible human review rather than through a fixed score set in advance.

This dynamic behavior also appeared in a prior user study of distributed AI in SAR teams, in which the PerceptiSync framework was deployed within a scenario-driven, multi-role simulation that compared AI-assisted and non-assisted conditions [18]. Participants selected trust configurations before the study, including trust level, historical accumulation, and minimum trusted events, which directly influenced how reliability scores were assessed between users. Quantitative results showed that more moderate or trusting configurations were associated with higher Average Trust Scores (ATS) and shorter decision times, whereas more cautious configurations were associated with lower trust scores and longer deliberation. Qualitative feedback further suggested that trust assessments sometimes served as shared reference points that shaped team dialogue and coordination, while other teams relied more heavily on personal judgment.

Beyond user preferences, PerceptiSync prioritizes human judgment through an optional *trust monitoring* feature that enables users to inspect received information and override system-computed reliability values when necessary. This becomes especially important during breakdowns, such as when distributed sources disagree, a detection is weakly supported, or operational context suggests that the AI is missing something important. In these moments, HCSR should function not only as a descriptive score but also as a trigger for oversight. When enabled, users are presented with visualizations of both their perceived scene and information from the sending source, along with the system's assessed reliability. Users may then adjust or replace the computed value based on their real-time interpretation. While this can introduce bias, it enables subject matter experts (SMEs) such as SAR team leaders, drone operators, or canine handlers to correct cases where automated reliability assessments are wrong or incomplete due to occlusions, sensor blind spots, or operational context not represented in the model [24]. More broadly, SME oversight can also help regulate what information is shared, and under what constraints, balancing openness with security and privacy in high-stakes collaboration [25]. Framing resilience around human-centered oversight helps reposition breakdowns as manageable moments for intervention rather than as failures of the system alone.

4. Shift 3: Situated Explainability in XR

Explainable AI (XAI) has traditionally been delivered through post-hoc explanations such as textual descriptions, saliency maps, or dashboard-style interfaces [26]. While valuable for model debugging and offline inspection, these forms of explanation are often too demanding for embodied, time-critical work [27]. In high-risk settings, professionals cannot divert attention to interpret lengthy explanations. They need information that can be perceived at a glance and interpreted in context [14]. This points to a need for XAI as an integrated, perceptual cue embedded directly within the XR interface. In this paper, such cues are most useful when they help users recognize that reliability has changed, identify which source or detection is affected, and determine whether intervention is needed during a breakdown.

Mixed reality offers a medium for such situated explanations by enabling digital information to be spatially anchored to physical objects, locations, and teammates. Prior work shows that spatial overlays can support rapid sensemaking by presenting relevant cues directly in a user's field of view. As shown

in our XR user studies for challenging environments, KHAIT leverages first-person canine video streams as XR overlays to enable handlers to interpret hazards and detect survivors in real time [28], while Ajna introduces shared perception overlays that fuse detections from multiple sensing modalities into a common spatial frame for collaborative interpretation [1]. Building on these examples, we argue that reliability cues can serve as explanations in their own right. Unlike post-hoc explanations, which mainly clarify why an output was produced, situated explanations in XR can communicate how that output should be interpreted in context and whether it should be acted on, treated cautiously, or reviewed further. In this sense, an HCSR cue can serve as a compact explanation of current system reliability at the moment of use.

Our prior SAR user study using PerceptiSync trust cues illustrates both the promise and the present limitations of this idea [18]. In that study, trust assessments supported collaboration and selective reliance across distributed AI agents, but they were conveyed primarily through verbal or textual forms rather than embodied XR cues. This leaves an important design question: how should reliability be visualized to support resilience during a system breakdown? We argue that situated explainability is the interface-level mechanism that enables this. When incoming information conflicts with a user's local view, when multiple sources disagree, or when evidence remains too uncertain to justify autonomous acceptance, the XR interface can surface that state directly at the point of use. Rather than requiring users to consult a dashboard or infer failure indirectly, the system could attach glanceable cues to detections or sources that communicate degraded reliability, ambiguity, or a need for confirmation. These cues may take the form of peripheral highlights, small anchored widgets, alert-style overlays, or other lightweight spatial signals that indicate when verification, override, or delayed action may be warranted. Viewed this way, situated explainability in XR is not limited to revealing why AI produced an output; it also explains how that output should be interpreted and acted upon in the moment.

5. Implications and Open Challenges

One implication of this work is the need for more systematic methods to prototype, compare, and evaluate how HCSR operates within distributed XR systems before testing in challenging environments. User studies that evaluate AI systems and interfaces together can blur the distinction between system behavior and interface design, making it difficult to determine how reliability cues may influence trust calibration, intervention timing, and team coordination [29, 30]. In particular, evaluation should isolate the progression outlined in this paper: how evidence updates HCSR, how HCSR informs oversight during uncertainty or disagreement, and how situated cues shape user response. One promising way to fill this gap is the Council of Wizards (CoW), a novel Wizard-of-Oz method we previously proposed for studying distributed AI [18]. By pairing multiple human operators with Wizards to simulate distributed perception systems and information exchange, the CoW supports controlled study of how different reliability cues and oversight strategies shape sensemaking, coordination, and decision-making. This approach allows researchers to repeatedly test alternative designs for communicating HCSR, examine unintuitive visual cues, and study how users interpret and act on reliability cues under task demands.

How to visualize HCSR in XR so that it supports situated explainability remains an open challenge. Because these cues are intended to operationalize resilience during breakdowns, they should help users notice degraded reliability, localize the affected source or detection, and understand when confirmation or override is warranted. Reliability cues could be visualized through changes in color, audio, iconography, or subtle spatial objects embedded within the user's field of view. For example, a peripheral border or halo around the display could shift in color or intensity to signal cautious, moderate, or trusting reliability when viewing information originating from another user or AI system. Alternatively, reliability may be conveyed through small, glanceable XR widgets or alert-style overlays anchored to incoming detections, offering controls for validation or override. Determining which combinations of cues are most interpretable, least disruptive, and most supportive of HITL override remains an open research question.

6. Declaration on Generative AI

During the preparation of this work, the authors used X-GPT-5.2 for grammar and spell check. Further, the authors used X-AI-IMG for figure 1 to generate the image. After using these tools, the authors reviewed and edited the content as needed and takes full responsibility for the publication's content.

References

- [1] M. Wilchek, K. Luther, F. A. Batarseh, Ajna: A wearable shared perception system for extreme sensemaking, *ACM Trans. Interact. Intell. Syst.* 15 (2025). URL: <https://doi.org/10.1145/3690829>. doi:10.1145/3690829.
- [2] M. A. Livingston, L. J. Rosenblum, D. G. Brown, G. S. Schmidt, S. J. Julier, Y. Baillot, J. E. Swan, Z. Ai, P. Maassel, *Military Applications of Augmented Reality*, Springer New York, New York, NY, 2011, pp. 671–706. URL: https://doi.org/10.1007/978-1-4614-0064-6_31. doi:10.1007/978-1-4614-0064-6_31.
- [3] M. A. Livingston, Z. Ai, K. Karsch, G. O. Gibson, User interface design for military ar applications, *Virtual Reality* 15 (2011) 175–184. URL: <https://doi.org/10.1007/s10055-010-0179-1>. doi:10.1007/s10055-010-0179-1.
- [4] F. A. Batarseh, L. Freeman, C.-H. Huang, A survey on artificial intelligence assurance, *Journal of Big Data* 8 (2021) 60. URL: <https://doi.org/10.1186/s40537-021-00445-7>. doi:10.1186/s40537-021-00445-7.
- [5] G. C. Moore, I. Benbasat, Development of an instrument to measure the perceptions of adopting an information technology innovation, *Info. Sys. Research* 2 (1991) 192–222. URL: <https://doi.org/10.1287/isre.2.3.192>. doi:10.1287/isre.2.3.192.
- [6] M. Madsen, S. D. Gregor, Measuring human-computer trust, in: *Proceedings of the 11th Australasian Conference on Information Systems (ACIS 2000)*, Information Systems Management Research Centre, Queensland University of Technology, Brisbane, QLD, Australia, 2000, pp. 6–8. URL: <https://api.semanticscholar.org/CorpusID:18821611>.
- [7] O. Vereschak, G. Bailly, B. Caramiaux, How to evaluate trust in ai-assisted decision making? a survey of empirical methodologies 5 (2021). URL: <https://doi.org/10.1145/3476068>. doi:10.1145/3476068.
- [8] F. Ekman, M. Johansson, J. Sochor, Creating appropriate trust in automated vehicle systems: A framework for hmi design, *IEEE Transactions on Human-Machine Systems* 48 (2018) 95–101. doi:10.1109/THMS.2017.2776209.
- [9] M. Billinghurst, A. Clark, G. Lee, A survey of augmented reality, *Foundations and Trends in Human-Computer Interaction* 8 (2015) 73–272. URL: <https://doi.org/10.1561/1100000049>. doi:10.1561/1100000049. arXiv:<https://www.emerald.com/fthci/article-pdf/8/2-3/73/11031378/1100000049en.pdf>.
- [10] M. Billinghurst, H. Kato, Collaborative augmented reality, *Commun. ACM* 45 (2002) 64–70. URL: <https://doi.org/10.1145/514236.514265>. doi:10.1145/514236.514265.
- [11] K. A. Hoff, M. Bashir, Trust in automation: Integrating empirical evidence on factors that influence trust, *Human Factors* 57 (2015) 407–434. URL: <https://doi.org/10.1177/0018720814547570>. doi:10.1177/0018720814547570, PMID: 25875432.
- [12] J. D. Lee, K. A. See, Trust in automation: Designing for appropriate reliance, *Human Factors* 46 (2004) 50–80. URL: https://journals.sagepub.com/doi/abs/10.1518/hfes.46.1.50_30392. doi:10.1518/hfes.46.1.50_30392, PMID: 15151155.
- [13] S. M. Merritt, H. Heimbaugh, J. LaChapell, D. Lee, I trust it, but i don't know why: Effects of implicit attitudes toward automation on trust in an automated system, *Human Factors* 55 (2013) 520–534. URL: <https://doi.org/10.1177/0018720812465081>. doi:10.1177/0018720812465081, PMID: 23829027.
- [14] M. T. Dzindolet, S. A. Peterson, R. A. Pomranky, L. G. Pierce, H. P. Beck, The role of trust in

- automation reliance, *International Journal of Human-Computer Studies* 58 (2003) 697–718. URL: <https://www.sciencedirect.com/science/article/pii/S1071581903000387>. doi:[https://doi.org/10.1016/S1071-5819\(03\)00038-7](https://doi.org/10.1016/S1071-5819(03)00038-7), trust and Technology.
- [15] P. A. Hancock, D. R. Billings, K. E. Schaefer, J. Y. C. Chen, E. J. de Visser, R. Parasuraman, A meta-analysis of factors affecting trust in human-robot interaction, *Human Factors* 53 (2011) 517–527. URL: <https://doi.org/10.1177/0018720811417254>. doi:10.1177/0018720811417254.
- [16] D. Gunning, Darpa’s explainable artificial intelligence (xai) program, in: *Proceedings of the 24th International Conference on Intelligent User Interfaces, IUI ’19*, Association for Computing Machinery, New York, NY, USA, 2019, p. ii. URL: <https://doi.org/10.1145/3301275.3308446>. doi:10.1145/3301275.3308446.
- [17] M. Wilchek, M. Nguyen, Y. Wang, K. Luther, F. A. Batarseh, Perceptisync: Trustworthy object detection using crowds-in-the-loop for cyber-physical systems, *ACM Transactions on Cyber-Physical Systems* (2025). URL: <https://doi.org/10.1145/3746644>. doi:10.1145/3746644, just Accepted (advance online publication); volume/issue not yet assigned.
- [18] M. Wilchek, S. Dickinson, K. Luther, F. A. Batarseh, The influence of distributed ai in trust and collaboration for search-and-rescue teams, in: *Proceedings of the 2026 CHI Conference on Human Factors in Computing Systems (CHI ’26)*, Association for Computing Machinery, New York, NY, USA, 2026, pp. 1–20. URL: <https://doi.org/10.1145/3772318.3791523>. doi:10.1145/3772318.3791523, just Accepted.
- [19] M. Wilchek, W. Hanley, J. Lim, K. Luther, F. A. Batarseh, Human-in-the-loop for computer vision assurance: A survey, *Engineering Applications of Artificial Intelligence* 123 (2023) 106376. URL: <https://www.sciencedirect.com/science/article/pii/S0952197623005602>. doi:<https://doi.org/10.1016/j.engappai.2023.106376>.
- [20] S. C. Kohn, E. J. de Visser, E. Wiese, Y.-C. Lee, T. H. Shaw, Measurement of trust in automation: A narrative review and reference guide, *Frontiers in Psychology* 12 (2021) 604977. URL: <https://doi.org/10.3389/fpsyg.2021.604977>. doi:10.3389/fpsyg.2021.604977.
- [21] J. Guo, Q. Yang, S. Fu, R. Boyles, S. Turner, K. Clarke, Towards trustworthy perception information sharing on connected and autonomous vehicles, in: *2020 International Conference on Connected and Autonomous Driving (MetroCAD)*, 2020, pp. 85–90. doi:10.1109/MetroCAD48866.2020.00021.
- [22] D. Cirne, V. Calambur, Fostering trust and quantifying value of ai and ml, in: K. Arai (Ed.), *Proceedings of the Future Technologies Conference (FTC) 2024*, Volume 2, Springer Nature Switzerland, Cham, 2024, pp. 586–603.
- [23] S. Amershi, D. Weld, M. Vorvoreanu, A. Fourney, B. Nushi, P. Collisson, J. Suh, S. Iqbal, P. N. Bennett, K. Inkpen, J. Teevan, R. Kikin-Gil, E. Horvitz, Guidelines for human-ai interaction, in: *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, CHI ’19*, Association for Computing Machinery, New York, NY, USA, 2019, p. 1–13. URL: <https://doi.org/10.1145/3290605.3300233>. doi:10.1145/3290605.3300233.
- [24] M. De-Arteaga, R. Fogliato, A. Chouldechova, A case for humans-in-the-loop: Decisions in the presence of erroneous algorithmic scores, in: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems, CHI ’20*, Association for Computing Machinery, New York, NY, USA, 2020, p. 1–12. URL: <https://doi.org/10.1145/3313831.3376638>. doi:10.1145/3313831.3376638.
- [25] S. Venkatagiri, A. Gautam, K. Luther, Crowdsolve: Managing tensions in an expert-led crowdsourced investigation, *Proc. ACM Hum.-Comput. Interact.* 5 (2021). URL: <https://doi.org/10.1145/3449192>. doi:10.1145/3449192.
- [26] T. Nguyen, A. Canossa, J. Zhu, How human-centered explainable ai interface are designed and evaluated: A systematic survey, 2024. URL: <https://arxiv.org/abs/2403.14496>. arXiv:2403.14496.
- [27] J. Kim, H. Maathuis, D. Sent, Human-centered evaluation of explainable ai applications: a systematic review, *Frontiers in Artificial Intelligence Volume 7 - 2024* (2024). URL: <https://www.frontiersin.org/journals/artificial-intelligence/articles/10.3389/frai.2024.1456486>. doi:10.3389/frai.2024.1456486.
- [28] M. Wilchek, L. Wang, S. Dickinson, E. Feuerbacher, K. Luther, F. A. Batarseh, Khait: K-9 han-

dlar artificial intelligence teaming for collaborative sensemaking, in: Proceedings of the 30th International Conference on Intelligent User Interfaces, IUI '25, Association for Computing Machinery, New York, NY, USA, 2025, p. 925–937. URL: <https://doi.org/10.1145/3708359.3712107>. doi:10.1145/3708359.3712107.

- [29] S. Naveed, G. Stevens, D. Robin-Kern, An overview of the empirical evaluation of explainable ai (xai): A comprehensive guideline for user-centered evaluation in xai, *Applied Sciences* 14 (2024). URL: <https://www.mdpi.com/2076-3417/14/23/11288>. doi:10.3390/app142311288.
- [30] A. Rechkemmer, M. Yin, When confidence meets accuracy: Exploring the effects of multiple performance indicators on trust in machine learning models, in: Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems, CHI '22, Association for Computing Machinery, New York, NY, USA, 2022. URL: <https://doi.org/10.1145/3491102.3501967>. doi:10.1145/3491102.3501967.