

Transparency of AI-XR Systems: Insights from Experts

Clara Maathuis
Open University of the Netherlands
The Netherlands
clara.maathuis@ou.nl

Dragos Datcu
Independent Research
The Netherlands
email@dragosdatcu.eu

Abstract—Recent advancements in the field of Extended Reality (XR) have increasingly integrated AI (Artificial Intelligence) technologies to enhance user experience and interactions. These are now crucial in XR applications in tasks like real-time object recognition, responsive environment generation, and adaptive intelligent behavior creation. Nevertheless, AI-XR systems need to adhere to social, legal, and ethical standards, norms, and values to ensure that their development, deployment, and use is done in a responsible and trustworthy manner. This assures building trust, preventing misusing sensitive information, and safeguarding user privacy and autonomy. To this end, transparency of AI-XR systems is fundamental for building and assuring trust. Accordingly, this research aims to understand and reflect on the perspectives that experts have in relation to building transparent AI-XR systems. To do so, a workshop with 14 field experts is conducted and the findings serve as a design framework with a rich palette of technical, socio-technical, and human-centered instruments that could be of use to researchers and practitioners in this domain.

Keywords—transparency, explainability, explainable AI, AI-XR, explainable AI-XR, transparent AI-XR.

I. INTRODUCTION

“XR allows us to create environments we’ve never experienced but have always imagined. It’s the playground of the mind.” (Douglas Trumbull)

Extended Reality (XR) is a rapidly advancing field that blends real and virtual elements to create immersive experiences. Ranging from fully immersive Virtual Reality (VR), where users are transported into entirely virtual worlds, to Augmented Reality (AR), which overlays digital components onto the real environment, XR technologies are pushing the boundaries of interaction and realism. Core features like immersion, interactivity, and high fidelity in graphics and sound combine to foster a powerful sense of presence, allowing users to feel part of these digital environments [1]. With successful applications in the educational domain that support student skill acquisition, engagement, and knowledge transfer [2], biomedicine and healthcare in tasks such as instrumentation assessment and diagnostic [3], assessment of cybersickness side effects and therapeutic applications that tackle symptoms such as disorientation, nausea, and oculomotor disturbance [4], XR systems are becoming more advanced and also more accessible

to the general public, with improvements in affordability, usability, and device portability.

Given recent technological developments, AI (Artificial Intelligence) techniques are increasingly used in various components of the XR environments, making them become more realistic and responsive. Nevertheless, while developing, deploying, and using AI-XR systems, ethical, social, and legal considerations are essential to ensure that they are responsible, trustworthy, and provide meaningful benefits to users and society [5]. An important value that AI-XR systems need to embed is transparency. This is essential for assuring users’ trust, engagement, and informed interaction within immersive environments, as they allow users to understand the rationale behind the decisions made and how these would impact their experience. Accordingly, transparent AI-XR systems need to provide explanations to their users in a human-centered approach, being clear and accessible, allowing them to maintain understanding, agency, address potential biases, foster safety, security, and privacy along their experience [6, 7].

While various recent studies and applications are dedicated to building transparent AI-XR systems, they are often tailored to a specific domain, context, user type, and less trying to grasp the multifaceted dimensions that characterize transparency and explainability of AI-XR systems from the perspective of field experts. Accordingly, this research aims to provide a multifaceted perspective of the meaning and implementation of transparency through explanations for AI-XR systems. It does that by means of conducting an extensive literature review and a workshop with fourteen field experts. In these lines, this research has the following two core contributions:

- To both research and practitioner communities engaged in building a comprehensive framework for transparency of AI-XR systems that captures the multifaceted and multidisciplinary aspects that characterize them in an adaptive, responsible, and human-centred manner.
- To decision-makers and designers involved in establishing and implementing policies and guidelines for developing and deploying transparent AI-XR systems through explainability of their decisions made across various applications, while addressing relevant complexities and challenges.

The remainder of the article is structured as follows. Section 2 discusses relevant studies to this domain. Section 3 presents the methodological approach considered to achieve the goals of this research. Section 4 discusses the findings obtained from the workshop conducted with field experts. At the end, concluding remarks are provided in Section 5.

II. RELATED RESEARCH

AI is known for its transformative power, bringing capabilities that surpass traditional boundaries and enabling novel interactive experiences in XR, including lifelike simulations, customized virtual agents, and real-time communication responsive learning environments [8, 9]. At the same time, applications of AI in the XR domain include its use for emotion recognition which enables XR systems to gauge user emotions through features like eye gaze and voice stress analysis, fostering more empathic interactions between digital avatars and users, and tracking user movements in real-time, which not only improves communication, but also creates a sense of co-presence among participants in a virtual environment [10]. Moreover, [11] argue that the synergy between AI and XR is characterized by four elements. First, AI confers intelligence to XR systems by enabling them to process and interpret complex data. Second, XR assists AI by providing immersive environments for training AI models, allowing for more effective learning and adaptation. Third, this integration facilitates building advanced visualization methods that enhance user experiences and decision-making processes across various domains, including medical diagnosis and military training.

In the context of capturing, representing, and embedding ethical values since the design of AI systems, transparency can be assured through the implementation of XAI methods that allow stakeholders and users to comprehend and trust the systems. This assures that AI systems are not perceived as ‘black boxes’ which can lead to skepticism and reluctance to adopt these technologies. Here the relationship between transparency and explainability is important: while explainability aims to make AI decisions understandable, transparency provides the necessary context and rationale behind those decisions [12]. Moreover, XAI methods can be categorized into three primary types: ante-hoc versus post-hoc, model-based versus model-agnostic, and global versus local explanations. Ante-hoc explanations involve models inherently designed for interpretability, such as decision trees, while post-hoc explanations, like LIME (Local Interpretable Model-agnostic Explanations) or SHAP (Shapley Additive exPlanations), are applied to interpret complex models after predictions are made. Model-based approaches are customized to specific model types, emphasizing interpretable structures, whereas model-agnostic methods can be applied universally across different model architectures. And global explanations provide insights into a model’s overall behavior, while local explanations examine individual predictions, clarifying specific decisions on a case-by-case basis [13]. The explanations provided need to be clear, ensuring they are

accessible to users of various technical backgrounds, and be meaningful, offering relevant insights that help users understand model’s decisions. Further, the explanations need to be accurate, reflecting model’s behavior, and highlight the model’s knowledge limits to inform users of contexts where it may perform poorly. Additionally, contextual relevance and comprehensiveness are also essential, ensuring explanations align with user needs and cover key aspects of the model’s decision-making process [14].

Transparent AI solutions are essential in the XR domain to build user trust, ensure safety, and address safety and privacy concerns. This helps users feel secure with systems that collect and process personal data, allowing them to understand how decisions are made [15], which is crucial in data-intensive XR environments like the metaverse. By making AI-XR systems interpretable and generalizable across various applications, developers can not only safeguard against security and privacy vulnerabilities and unethical practices, but also enhance user engagement and trust in these immersive systems [16]. Applied in healthcare, transparent AI-XR systems clarify the medical recommendations made, building trust for diagnostic applications. In education, they explain personalized learning paths, helping students understand their progress, while in gaming, it enhances immersion by elucidating AI opponent strategies. For virtual real estate, market trends and property values could be easier interpreted and further aiding buyer decisions. In virtual workspaces, transparent AI-XR systems promote team collaboration by explaining the task suggestions and enriching the user experience [17].

[18] explore the potential of XR as an enabler for XAI, developing a novel solution that receives positive user feedback and highlighted the effectiveness of XR for explaining AI systems. The intersection of XAI and XR is particularly impactful in specific application domains, such as healthcare, where [13] propose an architecture combining AI and blockchain in metaverse environments, using XAI techniques like GradCAM and LIME to provide logical reasoning for disease predictions while ensuring transparency.

The importance of transparency in AI-XR systems extends to user interaction and decision-making contexts. [19] explore AR explanation visualization designs, examining different approaches to feature importance analysis visualization in context-aware AR systems. Their findings provide valuable insights into the design of AR explanation artifacts that enhance user understanding of AI recommendations. Recent advancements in emotion detection within VR environments also emphasize the role of transparency. To this end, [20] use XAI techniques, specifically SHAP and LIME, to elucidate relationships between eye-tracking metrics and emotional states in VR, demonstrating how explainable approaches can bridge the gap between complex machine-learning models and practical applications. [21] propose the VR-LENS framework, emphasizing the importance of explainable approaches in cybersickness detection, while maintaining computational efficiency for deployment on standalone VR devices. Similarly, [22] explore the role of XAI in creating

personalized psychological, and emotional responses in VR environments through the Extended-Personal Reality concept.

Recent practitioner studies establish concrete standards for implementing transparent AI in XR systems. [18] conduct practical evaluations of an XAI-XR solution, providing industry-validated recommendations based on user responses. Their research sets baseline standards for practitioners implementing explainable AI in XR environments. The XAIR framework, developed through extensive practitioner research by [23], offers practical guidelines derived from large-scale user studies and expert workshops, particularly valuable for commercial deployments. This practitioner-focused approach is complemented by [24], who provides industry-validated frameworks for implementing algorithmic transparency in interactive systems. [25] proposes empirically validated approaches for implementing transparent AI assistants in VR environments for public transport applications. [26] extends these standards to biometric identification systems, establishing practical guidelines for transparent user identification in VR.

[22] contributed practical standards for implementing transparent personalization in commercial VR applications through their EXPERIENCE project. The development of future practice standards is being actively shaped through industry collaboration. [27] facilitated practitioner workshops to establish implementation guidelines for transparent AI-XR systems. These efforts are complemented by [28] who provide practical standards for implementing transparent behavioral analysis in virtual environments.

As this extensive literature review conducted shows, various research and practitioner efforts have been dedicated to understanding and implementing transparent AI-XR systems. Nevertheless, only a limited number of studies have been directly accounting for the direct integration of experts' perspective in their approach.

III. RESEARCH METHODOLOGY

To explore the development of transparent and explainable AI-XR systems, this research adopts a dual approach. On the one hand, it conducts an extensive literature review in this domain. On the other hand, it conducts a workshop with experts in this field in order to capture the perspectives that practitioners have in relation to this topic [29, 30]. The workshop was conducted in September 2024 being collocated with the ICUSI (International Conference on User-System Interaction) Conference. The participants of the conference were invited to engage and fourteen experts expressed their interest to participate and actively participated in the workshop. In this process, the following research questions (RQ) were posed and tackled:

RQ`1: How can AI-XR systems effectively communicate their decision-making processes to users in real-time without disrupting the immersive experience?

RQ2: What level of explanations should be provided to users regarding system decisions and actions?

RQ3: What level of technical detail should be provided in explanations of AI-XR systems operations to cater to users with varying levels of technical expertise?

RQ4: What role should user-controlled settings play in determining the level of explanations provided by AI-XR systems?

RQ5: What mechanism should be in place to allow users to query AI-XR systems about specific decisions or actions in the XR environment?

RQ6: What metrics or evaluation methods can be used to assess the effectiveness and user satisfaction with explanations provided by the AI-XR systems?

While the literature review component is conducted for synthesizing existing knowledge about building and enhancing transparency through explainability for AI-XR systems, the workshop reveals practical attributes, methods, and techniques that portray the taxonomy of this approach in an actionable manner. The workshop is conducted in English by the authors who are also engaged in the process of collecting and analysing the data. For each of the questions discussed, groups of answers are gathered and further clustered based on their meaning and relation to a specific category. Hence, the findings obtained are further presented in this article and discussed in relation to the existing state-of-the-art perspectives and applications developed in this domain of interest.

IV. TRANSPARENCY AND EXPLAINABILITY

In relation to the communication dimension (RQ1), the results suggest that AI-XR systems should communicate their decisions in a transparent manner without compromising immersion, emphasizing the need for multimodal and user-centred feedback within XR environments. Key techniques include using virtual agents, sound-based cues, haptic feedback, and visual indicators, each aiming to integrate decision transparency seamlessly into the user experience. For example, real-time communication via virtual agents or chatbots is suitable for providing natural and immediate feedback within the environment, helping users to interpret the choices and decisions that AI makes as part of the narrative. Haptic feedback can subtly show the importance without requiring the user's visual or cognitive attention. Moreover, volume and tone changes in environmental sounds could be used as non-verbal cues to capture the user attention in a gentle way. Nevertheless, these approaches need to be carefully calibrated to avoid sensory overload, as excessive or misaligned signals could detract from the immersive qualities of the XR experience. At the same time, decision-making transparency could also be enhanced by offering users customization options, such as choosing between auditory or visual feedback according to personal preference or task relevance. These aspects are illustrated in Figure 1 and show the importance of adaptivity and consideration of a transparent user-centred design which integrates both the environmental context and the user's cognitive load, ensuring the system provides relevant transparent feedback without overwhelming the sensory experience.

environment. In particular, voice commands are seen as effective for both simple and complex queries, offering a hands-free, immediate way to communicate with the system. Furthermore, gestural-based interaction that include hand, eye, or sign language would allow users to inquire the system to explain its decisions in ways that mirror real-world interactions. In addition to these methods, video-based predictions of actions and consequences are suggested to provide contextual insight into decisions. This functionality facilitates users to have a visual or interactive preview of potential outcomes based on AI actions, providing an intuitive explanation tool within the XR environment. Other recommendations, such as text-based interfaces for basic tasks or XR menus, could be tailored to varied use cases and complexity levels, balancing user needs for simplicity or depth. Moreover, the inclusion of Natural Language Processing (NLP) techniques further shows the role of conversational interfaces that are adaptive to user intent and language, making the system's responses more personalized. Therefore, these mechanisms aim to enhance or foster transparency and engagement by allowing users to query the AI-XR system with minimal cognitive effort, supporting both the usability and trust in decisions made by these systems.

Figure 5. Query mechanisms Word Cloud

effectiveness; higher engagement and lower error rates can suggest that explanations are fostering a positive, efficient user experience. Together, these metrics create a holistic evaluation framework that balances the need for measurable task performance with user-centered feedback, ensuring that explanations meet functional requirements and enhance the user experience and interaction.

Figure 6. Evaluation methods and metrics Word Cloud

V. DISCUSSION AND CONCLUSIONS

To this end, this research aims to critically analyze existing trends and applications of transparent AI-XR systems in various domains, and to couple the essential elements with the perspective of field experts in this domain. The findings suggest a clear need for adaptable and user-centered design elements captured in the AI-XR systems that balance explanatory depth with immersion. Moreover, the experts emphasize the need for customizable explanations, recommending settings that allow them to select the detail, and format of AI decision rationales based on their expertise level and task complexity. At the same time, mechanisms for querying AI actions and their corresponding explanations in order to ensure intuitive, non-intrusive interactions within the XR environment. Accordingly, the evaluation of explanation effectiveness needs to incorporate both performance-based metrics (e.g., task completion time, accuracy) and subjective user feedback (e.g., satisfaction ratings, qualitative interviews) to capture how well explanations support understanding and efficiency. Therefore, an adaptive and responsive stance is necessary when developing transparent AI-XR systems, in this

way contributing to building a responsible, safe, and trustworthy environment.

REFERENCES

- [1] H. Lee, “A conceptual model of immersive experience in extended reality”, 2020.
- [2] N. Pellas, I. Kazanidis, and G. Palaigeorgiou, “A systematic literature review of mixed reality environments in K-12 education”, *Education and Information Technologies*, “vol. 25, no. 4, pp. 2481-2520, 2020.
- [3] J. Yuan, S. S. Hassan, J. Wu, C. R. Koger, R. R. S. Packard, F. Shi., ... and Y. Ding, Y, “Extended reality for biomedicine”, *Nature Reviews Methods Primers*, vol. 3, no. 1, pp. 14, 2023.
- [4] L. Simón-Vicente, S. Rodríguez-Cano, V. Delgado-Benito, V., Ausín-Villaverde, and E. C. Delgado, “Cybersickness. A systematic literature review of adverse effects related to virtual reality”, *Neurología (English Edition)*, pp. 39, no. 8, pp. 701-709, 2024.
- [5] T. Hirzle, F. Müller, F. Draxler, M. Schmitz, P. Knierim, and K. Hornbæk, “When xr and ai meet-a scoping review on extended reality and artificial intelligence”, in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pp. 1-45, 2023.
- [6] C. Norval, R. Cloete, and J. Singh, “Navigating the audit landscape: A framework for developing transparent and auditable XR”, in *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1418-1431, 2023.
- [7] C. Maathuis, “Human Centered Explainable AI Framework for Military Cyber Operations”, in *IEEE Military Communications Conference 2023*, pp. 260-267, IEEE, 2023.
- [8] V. D. Păvăloaia and S. C. Necula, “Artificial intelligence as a disruptive technology—a systematic literature review”, *Electronics*, vol. 12, no. 5, pp. 1102, 2023.
- [9] T. Huynh-The, Q. V. Pham, X. Q. Pham, T. T. Nguyen, Z. Han, and D. S. Kim, “Artificial intelligence for the metaverse: A survey”, *Engineering Applications of Artificial Intelligence*, no. 117, pp. 105581, 2023.
- [10] A. Grech, J. Mehne and A. Wodehouse, “An extended AI-experience: Industry 5.0 in creative product innovation”, *Sensors*, vol. 23, no. 6, pp. 3009, 2023.
- [11] D. Reiners, M. R. Davahli, W. Karwowski and C. Cruz-Neira, “The combination of artificial intelligence and extended reality: A systematic review”, *Frontiers in Virtual Reality*, no. 2, pp. 721933, 2021.
- [12] L. Longo, M. Brcic, F. Cabitza, J. Choi, R. Confalonieri, J. Del Ser, ... and S. Stumpf, “Explainable Artificial Intelligence (XAI) 2.0: A manifesto of open challenges and interdisciplinary research directions”, *Information Fusion*, no. 106, pp. 102301, 2024.
- [13] S. Ali, T. Abuhmed, S. El-Sappagh, K. Muhammad, J. M. Alonso-Moral, R. Confalonieri, ... and F. Herrera, “Explainable Artificial Intelligence (XAI): What we know and what is left to attain Trustworthy Artificial Intelligence”, *Information fusion*, no. 99, pp. 101805, 2023.
- [14] A. Saranya and R. Subhashini, “A systematic review of Explainable Artificial Intelligence models and applications: Recent developments and future trends”, *Decision Analytics Journal*, no. 7, 2023.
- [15] C. Maathuis, M.A. Cidota, D. Datcu and L. Marin, “Explainable Artificial Intelligence Techniques for Extended Reality Systems: a Systematic Literature Review”, In: International Conference on User-System Interaction (ICUSI 2024), pp. 2-9, 2024.
- [16] A. Qayyum, M. A. Butt, H. Ali, M. Usman, O. Halabi, A. Al-Fuqaha, ... & J. Qadir, “Secure and trustworthy artificial intelligence-extended reality (AI-XR) for metaverses”, *ACM Computing Surveys*, vol. 56, no. 7, pp. 1-38, 2024.
- [17] G. Yenduri, G. Srivastava, M. Ramalingam, D. Reddy, M. Uzair and T. R. Gadekallu, “Explainable AI for the Metaverse: A Short Survey, in *International Conference on Intelligent Metaverse Technologies & Applications (iMETA)*, pp. 1-6, IEEE, 2023.
- [18] R. Wheeler and F. Carroll, “An Explainable AI Solution: Exploring Extended Reality As A Way To Make Artificial Intelligence More Transparent And Trustworthy”, in Onwubiko, C., et al. Proceedings of the International Conference on Cybersecurity, Situational Awareness and Social Media. Springer Proceedings in Complexity, 2023.
- [19] M. Zheng, R. J. Thomas, X. Pan, Z. Xu, Y. Liang, A.G. Campbell, “Augmenting Feature Importance Analysis: How Color and Size Can Affect Context-Aware AR Explanation Visualizations?”, in IEEE International Symposium on Mixed and Augmented Reality (ISMAR), 2022.
- [20] M. Bekler, M. Yilmaz, H.E. Ilgin, “Assessing Feature Importance in Eye-Tracking Data within Virtual Reality Using Explainable Artificial Intelligence Techniques”. *Appl. Sci.*, no. 14, 2024.
- [21] R.K. Kundu, R. Islam, P. Callyam, K.A. Hoque, “TruVR: Trustworthy Cybersickness Detection using Explainable Machine Learning”, 2022 IEEE International Symposium on Mixed and Augmented Reality (ISMAR), 2022.
- [22] G. Valenza, M. Alcañiz, V. Carli, G. Dudnik, C. Gentili, J.G. Provinciale, S. Rossi, N. Toschi and V. van Wassenhove, The EXPERIENCE Project: Automatic virtualization of "extended personal reality" through biomedical signal processing and explainable artificial intelligence [Applications Corner]. *IEEE Signal Processing Magazine*, vol. 41, no. 1, pp. 60-66., 2024.
- [23] X. Xu, M. Yu, T. Jonker, K. Todi, F. Lu, X. Qian, J. Belo, T. Wang, M. Li, A. Mun, T.Y. Wu, J. Shen, T. Zhang, N. Kokhlikyan, F. Wang, P. Sorenson, S.K. Kim, and H. Benko, “XAIR: A Framework of Explainable AI in Augmented Reality”. In: Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23), pp. 1–30, 2023.
- [24] T. Bitzer, M. Wiener and W. Cram, “Algorithmic Transparency: Concepts, Antecedents, and Consequences – A Review and Research Framework”. *Communications of the Association for Information Systems*, vol. 52, pp. 293-331, 2023.
- [25] A.K. Faulhaber, I. Ni and L. Schmidt, “The Effect of Explanations on Trust in an Assistance System for Public Transport Users and the Role of the Propensity to Trust”, In Proceedings of Mensch und Computer 2021 (MuC '21). Association for Computing Machinery, New York, NY, USA, pp. 303–310, 2021.
- [26] J. Liebers, P. Horn, C. Burschik, G. Uwe, S. Schneegass, “Using Gaze Behavior and Head Orientation for Implicit Identification in Virtual Reality”, 2021.
- [27] R. Suzuki, M. Gonzalez-Franco, M. Sra and D. Lindlbauer, “XR and AI: AI-Enabled Virtual, Augmented, and Mixed Reality”. In: Adjunct Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23 Adjunct), pp. 1–3, 2023.
- [28] A. Kalatian and B. Farooq, “Decoding pedestrian and automated vehicle interactions using immersive virtual reality and interpretable deep learning”, *Transportation Research Part C: Emerging Technologies*, vol. 124, 2021.
- [29] Miles, M. B., A. B. Huberman, and J. Saldaña, “*Qualitative Data Analysis: A Methods Sourcebook*”, 3rd ed. Los Angeles: Sage, 2014.
- [30] R. Ørngreen and K. T. Levinsen, “Workshops as a research methodology”, *Electronic Journal of E-learning*, vol. 15, no. 1, pp. 70-81, 2017.