
Towards Better User Requirements: How to Involve Human Participants in XAI Research

Thu Nguyen

IT University of Copenhagen
Denmark
irng@itu.dk

Jichen Zhu

IT University of Copenhagen
Denmark
jichen.zhu@gmail.com

Abstract

Human-Center eXplainable AI (HCXAI) literature identifies the need to address user needs. This paper examines how existing XAI research involves human users in designing and developing XAI systems and identifies limitations in current practices, especially regarding how researchers identify user requirements. Finally, we propose several suggestions on how to derive better user requirements.

1 User Requirements for XAI Explanation

The research field of eXplainable AI (XAI) has emerged to make AI more transparent and trustworthy to humans by opening the AI black-box and explaining its underlying operation[1]. While the field has made significant breakthroughs in *technical explainability*, it has limited success producing the *effective explanations* needed by users[2, 3]. As a result, most explanations produced by XAI still lack usability, practical interpretability, and efficacy for real users [4, 5, 6, 7, 8]. This viewpoint aligns with the one proposed by Liao and Kushney [3], where they argue that explanations should address stakeholders’ needs. In our paper, we focus on understanding explanation needs from *lay-end users*, who are direct users of the XAI system but have little domain and AI knowledge. This paper’s *argument* is that a key step towards human-centered XAI and effective explanations is better-defined user requirements through deeper engagement with human users to gather that information.

User requirements are insights from users about their needs, problems, and the context of use of an interactive system often derived from observing users performing tasks, interviewing users, and conducting focus groups [9]. It is an essential part of the established User-Centered Design (UCD) process for interactive systems in general. Recently there has been growing acknowledgment that user requirements for explanation are necessary to build a usable XAI system [3]. However, existing work typically attempted to define user requirements based on what the researchers identify as desired qualities of an XAI system [10, 11, 12]. While this approach helps to formalize the evaluation criteria, it risks bringing XAI researchers’ bias as it is not coming from end users[6]. Therefore, there is a need for methods to extract user requirements from end users directly.

According to the established UCD process, user requirements should be derived before the XAI system’s development to guide the development process towards an outcome that users will find valuable. However, in existing XAI research, only a small portion engages humans in deriving user requirements and other aspects of XAI development. Among them, many diverge from the established UCD design process regarding 1) when user requirements are collected, 2) how to collect user requirements for black-box models, and 3) which user groups are involved. The rest of this section will identify these issues, and the next section proposes approaches to address them.

With a few exceptions (e.g., [13, 14, 15, 16]), XAI development teams did not collect user requirements from human users at the early stage of XAI development yet. By contrast, a lot of existing research (e.g., [17, 18, 19, 20, 21, 22]) only engaged human users at the later stage to evaluate

the XAI systems. This practice means that a lot of current XAI research is driven by researchers' intuition, and users are only asked at the end to validate the design decisions.

Among these few projects that obtained user requirements, most seek to understand the user reasoning for Machine Learning models with high interpretability. For example, they focus on how users make sense of Naive Bayes [23, 22] and statistical modeling [16], and use that information to derive user requirements. By contrast, there is a limited understanding of users' sense-making of black box models such as Deep Neural Networks. This is problematic because, without such knowledge, it is difficult to identify effective ways to communicate the information from XAI to users in ways that fit their cognitive needs.

Finally, user requirements gathering in current XAI research typically does not target specific lay-end user groups. As user requirements of explanation highly depend on the context of use [9, 14], the user needs [3, 24], and user profile [9]. For example, AI developers need explanations to understand the inner working of the model to debug the algorithmic error [3]. In contrast, non-expert end users will be overwhelmed by the same amount of technical details. For the same user group, whether they use AI in a high-stake or low-stake context will also affect the type of explanation they need. Also, their ability to understand explanations depends on their knowledge of the application domain and AI. Nonetheless, there are many XAI research studies (e.g., [25, 26, 27, 21]) rely on Amazon Mechanical Turk (AMT) workers, who do not represent a specific domain, needs, or context of use. We argue that XAI user research should better consider the context of use, domain, and expertise of actual end users.

2 Proposed Approaches to Deriving Better HCXAI User Requirements

We first suggest that human-centered XAI research should involve human users at the early stage of development. While there is growing recognition of the importance of the context of use and user needs [3], a common practice in the field is to derive user requirements from predefined desiderata such as [28, 29]. For example, Lipton [28] proposed desiderata such as *trust*, *causality*, *transferability*, *informativeness*, and *fair and ethical decision making*. Bansal proposed *actionability* [29] as another desiderata for XAI. While these desiderata offer a general direction, they are often too generalized to target the specific requirements in a given context of use. For instance, what is considered trustworthy or informative can differ vastly from AI experts to novice end-users. Conducting user research, especially using qualitative methods (e.g., developing user group profiles, task models, and personas) at the early stage, can be a practical way to supplement the context-free desiderata.

Second, we argue that XAI user requirements should consider how to align users' mental models of the AI (i.e., how the user thinks the AI works) with the system's conceptual model (roughly speaking, how the AI works)[30]. Existing HCI research has established that the usability of a traditional system improves when users' mental models align with the system's conceptual model. Due to the low human interpretability of black-box systems, however, it is unclear whether users, especially non-technical users, are able to construct mental models that completely match the system's conceptual model[31]. There are currently no agreements over whether XAI should include all details of system logic[32, 33] or only selective important information[19, 34]). Researchers have identified cognitive science theories of how humans make sense of complex information as a potentially fruitful direction to design XAI systems[23]. We believe that more users studies of how users construct mental models of AI (e.g., [31]) and what cognitive-based design element (e.g., *affordance* [35]) can provide much-needed empirical evidence to bridge this knowledge gap.

Lastly, XAI researchers should identify specific target user groups and develop XAI systems that address their requirements. Any real-world XAI system will be used by specific people (e.g., AI experts, domain experts, non-experts), in particular contexts (work, play, health), and with specific needs. This is especially true because the explanation is social interaction and needs to adapt to the explainees[6]. To improve the usability and efficacy of XAI systems, researchers should identify the above key elements in user requirements. We believe that more research is needed in *explanation interfaces*[2, 24] to study how to effectively *communicate* the technical information XAI algorithms generate to users.

In conclusion, XAI research has yielded significant technical solutions to increase the interpretability and transparency of AI algorithms. To improve its real-world usefulness, we argue that more effort is needed to understand users of XAI systems, especially through better user requirements.

3 Acknowledgement

This work is supported by the Danish Novo Nordisk Foundation Grant NNF20OC0066119.

References

- [1] David Gunning. Explainable artificial intelligence (xai). *Defense advanced research projects agency (DARPA), nd Web*, 2(2):1, 2017.
- [2] Michael Chromik and Andreas Butz. Human-XAI Interaction: A Review and Design Principles for Explanation User Interfaces. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 12933 LNCS:619–640, 2021.
- [3] Q Vera Liao and Kush R Varshney. Human-centered explainable ai (xai): From algorithms to user experiences. *arXiv preprint arXiv:2110.10790*, 2021.
- [4] Ashraf Abdul, Jo Vermeulen, Danding Wang, Brian Y. Lim, and Mohan Kankanhalli. Trends and Trajectories for Explainable, Accountable and Intelligible Systems. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems - CHI '18*, pages 1–18, 2018.
- [5] Finale Doshi-Velez and Been Kim. A roadmap for a rigorous science of interpretability, 2017.
- [6] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence*, 267:1–38, 2019.
- [7] Jichen Zhu, Antonios Liapis, Sebastian Risi, Rafael Bidarra, and G. Michael Youngblood. Explainable AI for Designers: A Human-Centered Perspective on Mixed-Initiative Co-Creation. *IEEE Conference on Computational Intelligence and Games, CIG*, 2018-Augus, 2018.
- [8] Jichen Zhu, Jennifer Villareale, Nithesh Javvaji, Sebastian Risi, Mathias Löwe, Rush Weigelt, and Casper Hartevel. Player-ai interaction: What neural network games reveal about ai as play. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–17, 2021.
- [9] Thomas Geis, Knut Polkehn, Rolf Molich, and Oliver Kluge. Cpx-ur curriculum. *UXQB e. V*, 2016.
- [10] Nava Tintarev and Judith Masthoff. A survey of explanations in recommender systems. In *2007 IEEE 23rd international conference on data engineering workshop*, pages 801–810. IEEE, 2007.
- [11] Ingrid Nunes and Dietmar Jannach. A systematic review and taxonomy of explanations in decision support and recommender systems. *User Modeling and User-Adapted Interaction*, 27(3):393–444, 2017.
- [12] Oyindamola Williams. Towards human-centred explainable ai: A systematic literature review. 2021.
- [13] Upol Ehsan, Q Vera Liao, Michael Muller, Mark O Riedl, and Justin D Weisz. Expanding explainability: Towards social transparency in ai systems. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–19, 2021.
- [14] Q Vera Liao, Daniel Gruen, and Sarah Miller. Questioning the ai: informing design practices for explainable ai user experiences. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–15, 2020.
- [15] Vanessa Putnam and Cristina Conati. Exploring the need for explainable artificial intelligence (xai) in intelligent tutoring systems (its). In *IUI Workshops*, volume 19, 2019.
- [16] J Tullio. *How it works: A field study of non-technical users interacting with an intelligent system*. Hale, 2007.
- [17] Zhongxia Chen, Xiting Wang, Xing Xie, Mehul Parsana, Akshay Soni, Xiang Ao, and Enhong Chen. Towards explainable conversational recommendation. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 2994–3000, 2021.
- [18] Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and Payel Das. Explanations based on the missing: Towards contrastive explanations with pertinent negatives. *Advances in neural information processing systems*, 31, 2018.
- [19] Jonathan L. Herlocker, Joseph A. Konstan, and John Riedl. Explaining collaborative filtering recommendations. *Proceedings of the 2000 ACM conference on Computer supported cooperative work - CSCW '00*, pages 241–250, 2000.
- [20] Emilee Rader, Kelley Cotter, and Janghee Cho. Explanations as mechanisms for supporting algorithmic transparency. In *Proceedings of the 2018 CHI conference on human factors in computing systems*, pages 1–13, 2018.
- [21] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.

- [22] Simone Stumpf, Vidya Rajaram, Lida Li, Weng Keen Wong, Margaret Burnett, Thomas Dietterich, Erin Sullivan, and Jonathan Herlocker. Interacting meaningfully with machine learning systems: Three experiments. *International Journal of Human Computer Studies*, 67:639–662, 8 2009.
- [23] Danding Wang, Qian Yang, Ashraf Abdul, and Brian Y. Lim. Designing theory-driven user-centric explainable AI. In *CHI '19 Proceedings of the 2019 annual conference on Human factors in computing systems*, 2019.
- [24] Sina Mohseni, Niloofar Zarei, and Eric D Ragan. A multidisciplinary survey and framework for design and evaluation of explainable ai systems. *ACM Transactions on Interactive Intelligent Systems (TiIS)*, 11(3-4):1–45, 2021.
- [25] Hao-Fei Cheng, Ruotong Wang, Zheng Zhang, Fiona O’Connell, Terrance Gray, F Maxwell Harper, and Haiyi Zhu. Explaining decision-making algorithms through ui: Strategies to help non-expert stakeholders. In *Proceedings of the 2019 chi conference on human factors in computing systems*, pages 1–12, 2019.
- [26] Menaka Narayanan, Emily Chen, Jeffrey He, Been Kim, Sam Gershman, and Finale Doshi-Velez. *arXiv preprint arXiv:1802.00682*, 2018.
- [27] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [28] Zachary C Lipton. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018.
- [29] Gagan Bansal. Explanatory Dialogs: Towards Actionable, Interactive Explanations. In *AIES 2018 - Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 356–357, 2018.
- [30] Donald A Norman. *The psychology of everyday things*. Basic books, 1988.
- [31] Jennifer Villareale, Casper Hartevelde, and Jichen Zhu. "I Want To See How Smart This AI Really Is": Player Mental Model Development of an Adversarial AI Player. In *Proceedings of the 2022 ACM Conference of CHI PLAY*, 2022.
- [32] Todd Kulesza, Margaret Burnett, Weng-keen Wong, and Simone Stumpf. Principles of explanatory debugging to personalize interactive machine learning. In *Proceedings of the 20th international conference on intelligent user interfaces.*, pages 126–137, 2015.
- [33] Todd Kulesza, Simone Stumpf, Margaret Burnett, Sherry Yang, Irwin Kwan, and Weng Keen Wong. Too much, too little, or just right? Ways explanations impact end users’ mental models. *Proceedings of IEEE Symposium on Visual Languages and Human-Centric Computing, VL/HCC*, pages 3–10, 2013.
- [34] James Schaffer, Prasanna Giridhar, Debra Jones, Tobias Höllerer, Tarek Abdelzaher, and John O’donovan. Getting the message? A study of explanation interfaces for microblog data analysis. In *International Conference on Intelligent User Interfaces, Proceedings IUI*, volume 2015-Janua, pages 345–356, 2015.
- [35] Rex Hartson. Cognitive, physical, sensory, and functional affordances in interaction design. *Behaviour & information technology*, 22(5):315–338, 2003.