

Bidirectional Human-AI Alignment: Emerging Challenges and Opportunities

Hua Shen
University of Washington
Seattle, WA, USA
huashen@uw.edu

Tiffany Kneareem
Google
Cambridge, MA, USA
tkneareem@google.com

Reshmi Ghosh
Microsoft
Cambridge, MA, USA
reshmighosh@microsoft.com

Michael Xieyang Liu
Google DeepMind
Pittsburgh, PA, USA
lxieyang.ggl@gmail.com

Andrés Monroy-Hernández
Princeton University
Princeton, NJ, USA
andresmh@princeton.edu

Tongshuang Wu
Carnegie Mellon University
Pittsburgh, PA, USA
sherryw@cs.cmu.edu

Diyi Yang
Stanford University
Stanford, CA, USA
diyi@cs.stanford.edu

Yun Huang
University of Illinois at
Urbana-Champaign
Urbana-Champaign, IL, USA
yunhuang@illinois.edu

Tanushree Mitra
University of Washington
Seattle, WA, USA
tmitra@uw.edu

Yang Li
Google DeepMind
Mountain View, CA, USA
liyang@google.com

Marti A. Hearst
University of California, Berkeley
Berkeley, CA, USA
hearst@berkeley.edu

Abstract

Recent advancements in general-purpose AI have highlighted the urgent need to align AI systems with the goals, ethical principles, and values of individuals and society. Existing alignment research has been primarily approached as an AI-centered, static, and unidirectional process. However, this unidirectional perspective falls short of taking into account the dynamic and evolving interaction between humans and AI, necessitating a shift toward a bidirectional, interconnected mode of human-AI alignment. This SIG aims to outline the emerging areas of bidirectional human-AI alignment research, propose a blueprint of future goals and challenges for fundamental alignment research, and establish a shared platform to bring together experts from HCI, AI, social sciences, and more to advance interdisciplinary research and collaboration on human-AI alignment.

CCS Concepts

• **Human-centered computing** → **Interaction design**; **Human computer interaction (HCI)**; • **Social and professional topics** → **Computing / technology policy**; • **Security and privacy** → **Human and societal aspects of security and privacy**.

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

CHI '25, April 26 - May 01, 2025, Yokohama, Japan

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Keywords

bidirectional human-AI alignment, value-centered design, human-AI interactive alignment, social impact of AI alignment

ACM Reference Format:

Hua Shen, Tiffany Kneareem, Reshmi Ghosh, Michael Xieyang Liu, Andrés Monroy-Hernández, Tongshuang Wu, Diyi Yang, Yun Huang, Tanushree Mitra, Yang Li, and Marti A. Hearst. 2025. Bidirectional Human-AI Alignment: Emerging Challenges and Opportunities. In *Proceedings of The ACM CHI conference on Human Factors in Computing Systems (CHI '25)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

The rapid advancements in general-purpose AI has precipitated the urgent need to align these systems with values, ethical principles, and goals of individuals and society at large. This need, commonly referred to as “AI alignment,” [22, 24] is crucial for ensuring that AI systems function in a manner that is not only effective but also consistent with human values, minimizing harm and maximizing societal benefits. Traditionally, AI alignment has been viewed as a static, one-way process, with a primary focus on shaping AI systems to achieve desired outcomes and prevent negative side effects [5, 13, 15]. However, as AI systems become more integrated into everyday life and take on more complex decision-making roles, this unidirectional approach is proving inadequate [1]. AI systems interact with humans in evolving, unpredictable ways, generating feedback loops that influence both AI behavior and human responses. This **dynamic interaction necessitates a shift in how we think about alignment** [1]—one that recognizes the bidirectional and adaptive nature of human-AI relationships [18]. Rather than a one-time process or static goal, we should consider

alignment as a continuous, evolving engagement between humans and AI, requiring constant reassessment and recalibration. The concept of *Bidirectional Human-AI Alignment* offers a paradigm shift in how we approach the challenge of human-AI alignment [18], which emphasizes the **dynamic, mutual alignment process**. Particularly, it not only involves an AI-centered perspective, focuses on integrating human specifications into training, steering, and customizing AI. Also, it takes a human-centered perspective into account, aiming to preserve human agency and empower people to think critically when using AI, collaborate effectively with it, and adapt societal approaches to maximize its benefits for humanity. In this way, alignment becomes an interactive, reciprocal process. This ongoing dialogue between humans and machines is essential to achieving true alignment, as it allows both parties to evolve in response to changing contexts, goals, and ethical considerations. To work towards bidirectional human-AI alignment, grounded on our framework from a systematic survey of over 400 alignment papers [18], the core SIG objectives are twofold: (1) **broadening the current research of AI alignment** by expanding to adapting AI to dynamic human needs and societal contexts through research on value-sensitive design on AI, interactive alignment, alignment evaluation and social impacts, dynamic involvement for alignment, and more. (2) **fostering interdisciplinary collaboration** between researchers in multi-disciplinary domains, such as AI, HCI, and social sciences, creating a platform for exchange and innovation. Consequently, this SIG offers a comprehensive perspective and invites researchers to explore diverse aspects of alignment.

2 SIG Goals

The goal of this SIG includes outlining existing and emerging areas of bidirectional human-AI alignment, proposing a roadmap of future goals and challenges for alignment research. Specifically, this SIG primarily aims to explore four research questions:

- G.1 What are the fundamental individual and societal needs in bidirectional human-AI alignment?
- G.2 How can these alignment needs be addressed through the development of interaction and user experience mechanisms?
- G.3 How can we evaluate the effectiveness and societal impacts of bidirectional human-AI alignment?
- G.4 How will human-AI alignment co-evolve with dynamic changes, and what strategies can be designed to adapt to these changes?

Addressing these research questions requires a diverse group of interdisciplinary researchers and practitioners to work together. This SIG aims to bring together experts from HCI, AI, psychology, social sciences, and more domains to advance interdisciplinary research and collaboration on bidirectional human-AI alignment.

3 Themes

Given the four aforementioned SIG goals for exploring bidirectional human-AI alignment, we outline four corresponding themes that will be implemented in this SIG's events, including lightning talks, panel, group discussions, and post-SIG paper or blog outputs.

- **Alignment Objectives.** This theme explores the fundamental individual and societal needs in bidirectional human-AI alignment. Typical questions include: What human values should AI align with? How do humans' and AI's cognitive mental models

Time	Events
5 min	Welcome & Introduction
30 min	Lightning Talks + A Short Panel
5 min	Transition & Preparation
30 min	Theme-based Group Discussions
5 min	Summarization & Closing remarks

Table 1: Schedule for the 75-minute CHI Bi-Align SIG session.

differ during interactions? In what ways do individuals and societal groups vary in their approaches to aligning with AI?

- *Keywords & Example Papers:* value-centered alignment, alignment for individuals or societal groups, cognitive and behavioral alignment, alignment in multimodal AI, etc. [19]
- **Alignment Methods.** This theme investigates various approaches, such as interaction mechanisms and user experience design, to align AI with diverse human individuals and groups. Typical questions include: How can we interactively elicit human values for AI alignment? What interactive user interface designs can mitigate misalignment between humans and AI?
 - *Keywords & Example Papers:* interactive alignment, specification and requirements, UX design for AI alignment [23, 26].
- **Alignment Evaluation.** This theme examines methods for evaluating the effectiveness of human-AI alignment and its societal impacts across various applications. Key questions include: How can we quantitatively and qualitatively assess alignment between humans and AI? What are the impacts of AI alignment on individuals and societal groups?
 - *Keywords & Example Papers:* alignment evaluation, societal impacts, alignment for responsible and ethical AI, etc. [2]
- **Dynamic Alignment Evolvement.** This theme focuses on the dynamic co-evolution of human-AI alignment, addressing how alignment changes over time and how to adapt to these changes effectively.
 - *Keywords & Example Papers:* alignment evolution, adaptation in alignment, etc. [1]

4 Session Plan

4.1 Before the SIG

Before the SIG session, we will invite participants through social media promotion and professional mailing lists. We expect attendees from diverse backgrounds with different levels of familiarity and seniority with the topic. Furthermore, we have created a Slack Channel and a Google Folder for attendees to share materials and networking online.

4.2 SIG Schedule

We propose a 75-minute long session, with program consisting of following activities, in a hybrid format. While we encourage in-person attendance, synchronous online access will be provided. The tentative workshop schedule is detailed in Table 1.

Lightning Talks by Experts. We will invite four experts as our speakers to introduce the four SIG themes, and further provide a 5-min lightning talk tutorial for each theme. Particularly, we will ask the four speakers to collect representative papers in each theme and

present a short overview of each theme's research. The lightning talks last for 20 minutes in total.

A Short Panel with Speakers. The discussion panel will include the four experts with balanced perspectives from academia and industry who give the tutorials of four themes. The panel includes the Q&A among the moderator, panelists, and the audiences who attend the SIG in-person or online. The panel will last for 10 minutes.

Group Activity. We plan to host a 30-minute event of "On-the-spot Paper Writing" group activity. The groups will be formed using the "birds of a feather" format, which allows for individuals to join one of the four groups with the pre-defined themes. Our initial ideas of the group activity is:

- **On-the-spot Paper Writing:** Participants will choose from the predefined four alignment themes and join corresponding groups. Each group will brainstorm a research idea or an imaginary paper on their chosen topic. Deliverables may include an abstract, teaser figures illustrating key concepts, or compelling use cases. Groups will then share their work with others for feedback. At least one organizer will join each group for organizing the event and working on the imaginary paper.

4.3 Intended Community & Expected Size

We expect the SIG attendees to be academic or industrial researchers and practitioners who are broadly interested in alignment topics. Participants may come from a variety of disciplines, including human-computer interaction, Artificial Intelligence, machine learning, psychology, social science, and more. Deep technical expertise in AI is not a prerequisite. To bring interdisciplinary researchers together, we are also hosting a Bidirectional Human-AI Alignment Workshop accepted by the ICLR 2025 conference. Our shared Slack Channel enables researchers in CHI SIG, ICLR workshop, and more who are interested in alignment topics to easily connect and collaborate in the future. This SIG will particularly benefit those keen on exploring alignment-related topics including human-AI interaction, human-centered AI, responsible AI, AI for social good, user experience for AI, AI in various applications such as creativity, education, and more.

In order to facilitate meaningful, in-depth conversation, we tailor this SIG for a group of **30-50** participants, including 30-40 in-person and 10-20 remote participants. We may slightly adjust the SIG structure to accommodate a higher number of participants if needed.

5 Post-SIG Plans & Expected Outcomes

We plan to compile a comprehensive report summarizing the SIG's key discussions, presentations, and findings, which will be shared on open-access platforms such as ArXiv and our SIG website to reach a broad audience. Additionally, we aim to expand the outcomes of the *On-the-spot Paper Writing* session into full papers for submission to HCI conferences such as CHI. To facilitate ongoing communication and collaboration, we will establish dedicated Slack channels for the audience. Each thematic channel will be moderated by at least one co-organizer, who will act as a coordinator to guide discussions and activities during and after the SIG events.

References

- [1] Micah Carroll, Davis Foote, Anand Siththaranjan, Stuart Russell, and Anca Dragan. 2024. AI Alignment with Changing and Influenceable Reward Functions. *arXiv:2405.17713* (2024).
- [2] Preetam Prabhu Srikar Dammu, Hayoung Jung, Anjali Singh, Monojit Choudhury, and Tanushree Mitra. 2024. "They are uncultured": Unveiling Covert Harms and Social Threats in LLM Generated Conversations. *arXiv preprint arXiv:2405.05378* (2024).
- [3] Xuefeng Du, Reshmi Ghosh, Robert Sim, Ahmed Salem, Vitor Carvalho, Emily Lawton, Yixuan Li, and Jack W Stokes. 2024. VLMGuard: Defending VLMs against Malicious Prompts via Unlabeled Data. *arXiv preprint arXiv:2410.00296* (2024).
- [4] Peitong Duan, Bjoern Hartmann, Karina Nguyen, Yang Li, Marti Hearst, and Meredith Ringel Morris. 2023. Towards Semantically-Aware UI Design Tools: Design, Implementation, and Evaluation of Semantic Grouping Guidelines. (2023).
- [5] Nitesh Goyal, Minsuk Chang, and Michael Terry. 2024. Designing for Human-Agent Alignment: Understanding what humans want from their agents. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–6.
- [6] Hayoung Jung, Prerna Juneja, and Tanushree Mitra. 2024. Algorithmic Behaviors Across Regions: A Geolocation Audit of YouTube Search for COVID-19 Misinformation between the United States and South Africa. *arXiv preprint arXiv:2409.10168* (2024).
- [7] Philippe Laban, Jesse Vig, Marti A Hearst, Caiming Xiong, and Chien-Sheng Wu. 2023. Beyond the chat: Executable and verifiable text-editing with llms. *arXiv preprint arXiv:2309.15337* (2023).
- [8] Sangmin Lee, Minzhi Li, Bolin Lai, Wenqi Jia, Fiona Ryan, Xu Cao, Ozgur Kara, Bikram Boote, Weiyan Shi, Diyi Yang, et al. 2024. Towards Social AI: A Survey on Understanding Social Interactions. *arXiv preprint arXiv:2409.15316* (2024).
- [9] Gang Li and Yang Li. 2022. Spotlight: Mobile ui understanding using vision-language models with a focus. *arXiv preprint arXiv:2209.14927* (2022).
- [10] Michael Xieyang Liu, Tongshuang Wu, Tianying Chen, Franklin Mingzhe Li, Aniket Kittur, and Brad A Myers. 2024. Selenite: Scaffolding Online Sensemaking with Comprehensive Overviews Elicited from Large Language Models. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '24). Association for Computing Machinery, New York, NY, USA, Article 837, 26 pages. <https://doi.org/10.1145/3613904.3642149>
- [11] Kyle Lo, Joseph Chee Chang, Andrew Head, Jonathan Bragg, Amy X Zhang, Cassidy Trier, Chloe Anastasiades, Tal August, Russell Authur, Danielle Bragg, et al. 2023. The semantic reader project: Augmenting scholarly documents through ai-powered interactive reading interfaces. *arXiv preprint arXiv:2303.14334* (2023).
- [12] Qianou Ma, Hua Shen, Kenneth Koedinger, and Tongshuang Wu. 2024. How to Teach Programming in the AI Era? Using LLMs as a Teachable Agent for Debugging. *25th International Conference on Artificial Intelligence in Education (AIED 2024)* (2024).
- [13] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems* 35 (2022), 27730–27744.
- [14] Savvas Petridis, Michael Xieyang Liu, Alexander J Fiannaca, Vivian Tsai, Michael Terry, and Carrie J Cai. 2024. In Situ AI Prototyping: Infusing Multimodal Prompts into Mobile Settings with MobileMaker. *arXiv preprint arXiv:2405.03806* (2024).
- [15] Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect?. In *International Conference on Machine Learning*. PMLR, 29971–30004.
- [16] Hua Shen, Chieh-Yang Huang, Tongshuang Wu, and Ting-Hao Kenneth Huang. 2023. ConvXAI: Delivering heterogeneous AI explanations via conversations to support human-AI scientific writing. In *Companion Publication of the 2023 Conference on Computer Supported Cooperative Work and Social Computing*. 384–387.
- [17] Hua Shen and Ting-Hao Huang. 2020. How useful are the machine-generated interpretations to general users? a human evaluation on guessing the incorrectly predicted labels. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, Vol. 8. 168–172.
- [18] Hua Shen, Tiffany Kneare, Reshmi Ghosh, Kenan Alkiek, Kundan Krishna, Yachuan Liu, Ziqiao Ma, Savvas Petridis, Yi-Hao Peng, Li Qiwei, et al. 2024. Towards Bidirectional Human-AI Alignment: A Systematic Review for Clarifications, Framework, and Future Directions. *arXiv preprint arXiv:2406.09264* (2024).
- [19] Hua Shen, Tiffany Kneare, Reshmi Ghosh, Yu-Ju Yang, Tanushree Mitra, and Yun Huang. 2024. ValueCompass: A Framework of Fundamental Values for Human-AI Alignment. *arXiv preprint arXiv:2409.09586* (2024).
- [20] Hua Shen, Tongshuang Wu, Wenbo Guo, and Ting-Hao Huang. 2022. Are Shortest Rationales the Best Explanations for Human Understanding?. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. 10–19.

- [21] Jonathan Stray, Alon Halevy, Parisa Assar, Dylan Hadfield-Menell, Craig Boutilier, Amar Ashar, Chloe Bakalar, Lex Beattie, Michael Ekstrand, Claire Leibowicz, et al. 2024. Building human values into recommender systems: An interdisciplinary synthesis. *ACM Transactions on Recommender Systems* 2, 3 (2024), 1–57.
- [22] Michael Terry, Chinmay Kulkarni, Martin Wattenberg, Lucas Dixon, and Meredith Ringel Morris. 2023. AI Alignment in the Design of Interactive AI: Specification Alignment, Process Alignment, and Evaluation Support. *arXiv:2311.00710* (2023).
- [23] **Hua Shen**, Chieh-Yang Huang, Tongshuang Wu, and Ting-Hao 'Kenneth' Huang. 2023. ConvXAI: Delivering Heterogeneous AI Explanations via Conversations to Support Human-AI Scientific Writing. In *The 26th ACM Conference On Computer-Supported Cooperative Work And Social Computing - Demo (CSCW '23 Demo)*.
- [24] Wikipedia. 2024. AI alignment — Wikipedia, The Free Encyclopedia. <http://en.wikipedia.org/w/index.php?title=AI%20alignment&oldid=1220304776>. [Online; accessed 05-May-2024].
- [25] Robert Wolfe and Tanushree Mitra. 2024. The Impact and Opportunities of Generative AI in Fact-Checking. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 1531–1543.
- [26] Sherry Wu, Hua Shen, Daniel S Weld, Jeffrey Heer, and Marco Tulio Ribeiro. 2023. Scattershot: Interactive in-context example curation for text transformation. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*. 353–367.
- [27] Lei Zhang, Youjean Cho, Daekun Kim, Ava Robinson, Yu Jiang Tham, Rajan Vaish, and Andrés Monroy-Hernández. 2024. Interaction recording tools for creating interactive ar stories. US Patent App. 17/901,589.

A Supplementary Material

A.1 Organizers and Primary Contact

This workshop is organized by a team of experienced researchers with diverse expertise in human-AI alignment, representing various stages of their careers and spanning a wide range of demographic and geographical backgrounds. Collectively, the organizing committee has made significant contributions to both HCI and AI across a variety of topics, including alignment and values [18, 19, 21], AI explanation and sensemaking [10, 16, 17, 20], human-AI interaction [7, 11, 12, 14, 27], AI responsibility and auditing [3, 6], user experience design on AI and alignment [4, 9], and AI evaluation and social impact [8, 25]. The organizers have a strong track record of successfully hosting workshops, tutorials, and panels on human-centered AI, user experience and interaction design on AI, and more, bridging the HCI and AI communities.

Hua Shen (Primary Contact) is a postdoctoral scholar at the University of Washington. Her research centers on bidirectional human-AI alignment, aiming to empower humans to interactively explain, evaluate, and collaborate with AI, while incorporating human feedback and values to improve AI systems. She earned her Ph.D. from Pennsylvania State University and completed a postdoctoral fellowship at the University of Michigan. She initiated and led the systematic survey project of “Bidirectional Human-AI Alignment” [18].

Tiffany Kneareem (Co-Primary Contact) is a User Experience Researcher on the Material Design team at Google. Her research focus is on product designer-developer collaboration, creativity support tooling and opportunities for AI in the user interface (UI) design space. She holds a PhD in Information Sciences and Technologies from Pennsylvania State University, advised by Dr. John M. Carroll. She co-organized the CHI 2024 workshop on Computational UI.

Reshmi Ghosh is an Applied Scientist Lead for GenAI Safety in Microsoft’s Responsible AI and Security team. She was the core architect in LLM Safety, designing M365 CoPilots, and has previously worked on integrating machine learning features to Excel, Word, and PowerPoint. She graduated with a Ph.D. from Carnegie Mellon University, focusing on data reconstruction using NLP methods for mitigating climate change. She is a research advisor for teams at MIT CSAIL, UMass Amherst, UCLA, and Oxford University.

Michael Xieyang Liu is a research scientist at Google DeepMind. His research aims to improve human-AI interaction, with a particular focus on human interaction with multimodal large language models and controllable AI. Michael previously earned his Ph.D. from the Human-Computer Interaction Institute at Carnegie Mellon University, specializing in the intersection of HCI, programming tools, sensemaking, intelligent user interfaces, and human-AI interaction. Michael organized the Sensemaking workshop at CHI 2024.

Andrés Monroy-Hernández is an Assistant Professor co-leading the Princeton HCI Lab at Princeton University, where his research focuses on human-computer interaction and social computing. He is also an associated faculty at Princeton’s Center for Information Technology and Policy, the Keller Center for Innovation, the DeCenter, the Program in Cognitive Science, and the Program in Latin American Studies. Before Princeton, he led the HCI Research

team in the Microsoft Research’s FUSE Labs and Snap Research. He received his Ph.D. degree in Media Arts and Sciences from MIT.

Sherry Tongshuang Wu is an Assistant Professor in the Human-Computer Interaction Institute at Carnegie Mellon University. Her research lies at the intersection of Human-Computer Interaction and Natural Language Processing, aiming to design, evaluate, build, and interact with AI systems that are compatible with actual human goals. Before joining CMU, Sherry received her Ph.D. degree from the University of Washington. Sherry organized the TRAIT workshop at CHI 2022, 2023 and TREW workshop at CHI 2024.

Diyi Yang is an Assistant Professor in the Computer Science Department at Stanford University, affiliated with the Stanford NLP Group, Stanford HCI Group, Stanford AI Lab (SAIL), and Stanford Human-Centered Artificial Intelligence (HAI). Her research focuses on Socially Aware Natural Language Processing, aiming to better understand human communication in social context and build socially aware language technologies to support human-human and human-computer interaction. She received her Ph.D. degree in Language Technologies Institute at CMU.

Yun Huang is an Associate Professor in the School of Information Sciences at the University of Illinois at Urbana-Champaign. She is dedicated to innovating AI-based solutions that foster a synergistic relationship between humans and machines, enhancing educational opportunities to all and expanding access to community services. She received her Ph.D. in information and computer science from the University of California, Irvine.

Tanushree Mitra is an Associate Professor in the Information School at the University of Washington. Her research blends human-centered data science and social science principles to develop new knowledge, methods, and systems to defend against the epistemic risks of online mis(dis)information, bias, hate, and harm. She co-founded the Responsibility in AI Systems and Experiences (RAISE) Center at UW. She received her PhD in Computer Science from Georgia Tech’s School of Interactive Computing and her Masters in Computer Science from Texas A&M University.

Yang Li is a Senior Staff Research Scientist at Google DeepMind, and an affiliate faculty member at University of Washington. His research lies at the intersection of HCI and AI, focusing on general deep learning research and models for solving human interactive intelligence problems and improving user experiences. He earned a Ph.D. degree in Computer Science from the Chinese Academy of Sciences, and conducted postdoctoral research at UC Berkeley EECS. Yang organized multiple workshops that bridges HCI and AI/ML fields, including the first ICML AI&HCI workshop.

Marti A. Hearst is a Professor in the UC Berkeley School of Information and the Computer Science Division. She was Interim Dean and Head of School for the ISchool. Her research encompasses user interfaces with a focus on search, information visualization with a focus on text, computational linguistics, and educational technology. She co-founded the ACM Learning@Scale conference and is a former president of ACL, a CHI Academy member, the SIGIR Academy, an ACM Fellow, and an ACL Fellow. She received her PhD, MS, and BA degrees in Computer Science from UC Berkeley.

A.2 Organizing and Presenting Approach

We will host the SIG in a **hybrid format**, primarily in-person with an option for remote participation. All sessions will be live-streamed, with virtual breakout rooms available for remote attendees to join discussions. We will leverage the conference center’s standard equipment to meet the technical needs. A SIG website and Slack platform will serve as central hubs for engagement, offering details such as the program schedule, organizers, and speakers.

We provide asynchronous materials for all participants to access offline through both the SIG website and Slack platform. In case any technical or accessibility issues arise, we provide all important information, such as the program schedule, list of organizers and speakers, and pre-prints of accepted papers, on the SIG website. Besides, we allow all participants to engage in the SIG Slack for Q&A and discussion. Furthermore, we will release the videos of the SIG presentations on YouTube and list them on our SIG website.

A.3 Accessibility

We strive to create an inclusive workshop environment for all participants, including those with cognitive, mental health and physical disabilities. We will highly encourage authors to make their position paper accessible. For accepted papers, we plan to offer guidance on improving document accessibility, e.g., alt-text for images and tables, ensuring that the navigation hierarchy is intelligible for screen-readers. During the workshop, we will request that all participants follow accessibility best practices, for example use of a microphone at all times and turning on captioning for presentations.