# Human-AI Interaction Alignment: Designing, Evaluating, and Evolving Value-Centered AI For Reciprocal Human-AI Futures

Hua Shen
NYU Shanghai, New York University
China
huashen@nyu.edu

Tiffany Knearem
MBZUAI
United Arab Emirates
tiffany.knearem@mbzuai.ae.ac

Divy Thakkar
Google
United States
dthakkar@google.com

Pat Pataranutaporn
Massachusetts Institute of
Technologyy
United States
patpat@media.mit.edu

Anoop Sinha
Google, Paradigms of Intelligence
United States
anoopsinha@google.com

Yike Shi
Carnegie Mellon University, New
York University
United States
yikes@andrew.cmu.edu

Jenny Liang
Carnegie Mellon University
United States
jtliang@andrew.cmu.edu

Lama Ahmad
OpenAI
United States
lama@openai.com

Tanu Mitra
University of Washington
United States
tmitra@uw.edu

Brad A. Myers
Carnegie Mellon University
United States
bam@cs.cmu.edu

Yang Li
Google DeepMind
United States
liyang@google.com

## Abstract

The rapid integration of generative AI into everyday life underscores the need to move beyond unidirectional alignment models that only adapt AI to human values. This workshop focuses on bidirectional human-AI alignment, a dynamic, reciprocal process where humans and AI co-adapt through interaction, evaluation, and value-centered design. Building on our past CHI 2025 BiAlign SIG and ICLR 2025 Workshop, this workshop will bring together interdisciplinary researchers from HCI, AI, social sciences and more domains to advance value-centered AI and reciprocal human-AI collaboration. We focus on embedding human and societal values into alignment research, emphasizing not only steering AI toward human values but also enabling humans to critically engage with and evolve alongside AI systems. Through talks, interdisciplinary discussions, and collaborative activities, participants will explore methods for interactive alignment, frameworks for societal impact evaluation, and strategies for alignment in dynamic contexts. This workshop aims to bridge the disciplines' gaps and establish a shared agenda for responsible, reciprocal human-AI futures.

## CCS Concepts

• **Human-centered computing** → **HCI design and evaluation methods**; **Interaction paradigms**; **HCI theory, concepts and models**; **Interaction techniques**; **Interactive systems and tools**; **Empirical studies in HCI**.

## Keywords

bidirectional human-AI alignment, value-centered design, interactive alignment

## 1 Background and Motivation

Rapid advances in general-purpose and generative AI have precipitated the urgent need to align these systems with the values, ethical principles, and goals of individuals and society at large [14, 15, 19, 29, 32]. This need, commonly referred to as AI alignment [27, 30], is crucial to ensure that AI systems function in ways that are not only effective but also consistent with human values, minimizing harm and maximizing societal benefits [1, 8, 10, 16, 20, 21]. Traditionally, alignment has been viewed as a static, one-way process, with AI shaped to achieve desired outcomes and avoid negative side effects. Yet, as AI systems increasingly permeate daily life and take on complex decision-making roles, this unidirectional approach proves inadequate [2]. AI systems and humans interact in evolving and unpredictable ways, generating feedback loops that influence both AI behavior and human responses [5, 13, 17, 18].

This dynamic relationship necessitates a shift toward *bidirectional human-AI alignment* – a paradigm that treats alignment as a continuous, reciprocal process of mutual adaptation [23]. From this perspective, alignment requires not only steering AI toward human goals but also empowering humans to critically engage with, recalibrate, and evolve alongside AI systems.

**Past Workshops and Community Interest:** We have observed a growing number of workshops advancing AI alignment within the AI-centered research community, such as the NeurIPS 2024 Pluralistic Alignment Workshop, ICML 2024 Human Feedback for AI Alignment Workshop, and ICLR 2024 Representation Alignment Workshop. Yet, the **voice of the HCI community—bringing a human-centered perspective to alignment—has been notably absent**. To help address this gap, our team has begun building this vision through the ICLR 2025 Bidirectional Human-AI Alignment Workshop [26] and the CHI 2025 Bidirectional Human-AI Alignment SIG [24]. Both efforts attracted overwhelming interest: the ICLR workshop received more than 100 submissions, and the CHI SIG drew over 100 participants, filling the room to capacity. Many additional attendees who could not join requested access to materials, underscoring the strong demand for an interdisciplinary forum on this topic. The success of these past events in 2025 **highlights both the urgency of this research area and the opportunity for CHI to host a more expansive, dedicated workshop** where researchers can share insights, methods, and perspectives.

**Novel Contributions of CHI 2026 Workshop:** Building on this momentum, this CHI 2026 workshop – serving as the HCI home of our 2nd BiAlign Workshop – introduces several new initiatives to deepen engagement and expand the community: (1) *Interdisciplinary Integration*: Bridge interdisciplinary research via structured sessions and collaborations; (2) *Interactive Knowledge Creation*: Group activities such as collaborative paper prototyping, solution ideation, and concept mapping to co-develop new ideas; (3) *Expanded Accessibility*: A hybrid format supported by pre-recorded talks, shared materials, and a dedicated Slack channel to engage participants worldwide; (4) *Sustained Community Building*: Post-workshop initiatives including an open repository, ongoing discussion spaces, and opportunities for collective publications. Through these contributions, the workshop will advance value-centered approaches to bidirectional human-AI alignment, positioning CHI as a central venue for shaping reciprocal human-AI futures grounded in responsibility, values, and collaboration.

## 2 Workshop Goals and Themes

The overarching goal of this workshop is to establish a sustained, interdisciplinary forum for advancing bidirectional human-AI alignment — a paradigm that emphasizes dynamic, reciprocal processes of co-adaptation between humans and AI systems, grounded in human and societal values. Specifically, the workshop goals include:

G.1 **Operationalize Human and Societal Values:** We will identify and discuss frameworks for translating abstract values—such as fairness, agency, and responsibility—into actionable design principles and technical requirements. The workshop will surface practical strategies for embedding these values into the development and deployment of AI systems.

G.2 **Advance Design and Interaction Mechanisms:** Drawing from Human–Computer Interaction (HCI) methods, we will explore how interaction design, user experience research, and participatory approaches can enable humans to shape, critique, and recalibrate generative AI systems in real time. A particular emphasis will be placed on techniques that empower diverse stakeholders to meaningfully engage with, and guide, AI behavior.

G.3 **Explore Dynamic Evaluation Approaches:** The workshop will examine approaches to evaluating alignment at individual, community, and societal levels—balancing technical performance with societal impact. Participants will also consider strategies for sustaining alignment over time, recognizing that both humans and AI systems evolve as models acquire more advanced reasoning and adaptive capabilities.

G.4 **Foster Interdisciplinary Collaboration and Build Community:** Create a forum for researchers and practitioners in HCI, AI, and the social sciences to exchange perspectives, bridge disciplinary gaps, and shape shared research agendas. Through networking and collaborative activities, the workshop will strengthen the BiAlign community and lay the groundwork for sustained engagement beyond the event.

**Themes:** We will structure the workshop around four interrelated themes. Each theme will be introduced through short talks and exemplars, followed by interactive discussions and collaborative activities:

- **Value-Centered Alignment Objectives:** Explores which human and societal values should guide reciprocal human-AI alignment and how these values can be articulated and translated into technical and design processes.
  - *Research Questions:* What fundamental human and societal values should guide reciprocal human-AI alignment? In what ways might HCI contribute to the articulation and translation of values into technical and design processes?
  - *Keywords & Example Papers*: pluralistic values, human agency, cultural perspectives, value-sensitive design, etc [6, 22, 25].

- **Developing Interfaces and Interactions for Alignment:** Investigates design mechanisms—such as interfaces, interaction modalities, explanation systems, and participatory methods that empower humans to steer, critique, and co-create with AI systems.
  - *Research Questions*: What design mechanisms can help humans shape and steer AI behavior? What role do co-creation and participatory design methods play in aligning AI with evolving human needs? How do we uplift and retain human agency via effective human-AI collaboration?
  - *Keywords & Example Papers*: interactive alignment, UX for AI, participatory design, human-AI collaboration, etc [4, 7, 28, 31].

- **Evaluating Alignment and Societal Impacts:** Examines frameworks and methodologies for assessing bidirectional alignment, including both technical effectiveness and broader impacts such as trust and social well-being.
  - *Research Questions:* How should bidirectional alignment be measured—both technically and socially? What frameworks and methodologies can capture the broader impacts of alignment (e.g., trust, collective well-being, economic impact)?

- *Keywords & Example Papers*: alignment evaluation, societal impact, trust, responsible AI, etc [3, 10, 12].
- **Dynamic Co-Evolution of Human-AI Futures:** Considers alignment as an evolving process, highlighting strategies for sustaining reciprocal adaptation as both humans and AI change over time and across contexts.
  - *Research Questions*: How have alignment goals and practices evolved over time, as humans and AI systems mutually adapted? How can we envision and design for long-term reciprocal futures of human-AI collaboration?
  - *Keywords & Example Papers*: adaptability, resilience, lifelong learning, co-evolution, etc [2, 9, 11, 18].

Achieving these goals requires a diverse group of interdisciplinary researchers and practitioners, working together in open dialogue to shape, define and execute on alignment goals. This workshop 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 Organizers

This workshop brings together an interdisciplinary team of organizers spanning academia and industry with global representations from the United States, China, and United Arab Emirates. The organizers contribute expertise across HCI, CSCW, NLP, ML, psychology, AI safety and governance. The team includes researchers who have shaped the study of bidirectional human-AI alignment and working persistently to contribute this research to broader interdisciplinary communities. The organizers also bring substantial experience in organizing successful workshops and tutorials at premier venues such as CHI, CSCW, ICLR, EMNLP and more, ensuring effective facilitation and impactful outcomes.

**Hua Shen, Ph.D.** is an Assistant Professor of Computer Science at NYU Shanghai, New York University. Her work focuses on Bidirectional Human–AI Alignment: enabling humans to interactively explain, evaluate, and collaborate with AI systems, while integrating human feedback and values to improve AI systems. She has been recognized as a 2023 Rising Star in Data Science, and received multiple awards, including AIED 2024 Best Paper and Best Interactive Event, CSCW 2023 Best Demo, and IUI 2023 Best Paper Honorable Mention, and Google Research Scholarships. She is serving as Associate Chair for CHI 2026-2025, CHI LBW 2024, Program Committees for ACL, EMNLP, and more. She initiated the 2025 BiAlign CHI SIG and ICLR workshop, NeurIPS 2025 Human-AI Alignment Tutorial, CSCW 2025 Design for Hope workshop, EMNLP 2025 WiNLP workshop for HCI and AI communities.

**Tiffany Knearem, Ph.D.** is an Affiliated Assistant Professor at the Mohamed bin Zayed University of Artificial Intelligence (MBZUAI) and the head of TK Research, a UX and HCI research consultancy. She holds a Ph.D. in information sciences and technology from Pennsylvania State University, advised by Prof. John M. Carroll. Her recent research interests span human-AI alignment, AI-supported design workflows, and community informatics. She co-organized the CHI 2024 and CHI 2025 Computational UI workshops, and the CHI 2025 Bi-Align SIG.

**Divy Thakkar, Ph.D.** is a Staff Program Manager and Researcher at Google DeepMind, where he is building new interactions and

human-ai collaboration mechanisms for Gemini. His research has earned recognition at top HCI conferences, including CHI and CSCW. Thakkar completed his Ph.D. in Computer Science at City St. Georges, University of London.

**Pat Pataranutaporn, Ph.D.** is an Assistant Professor at the MIT Media Lab, where he directs the Cyborg Psychology Research Group and co-directs the Advancing Humans with AI (AHA) Program. His research develops AI systems that foster human flourishing, including pioneering studies on generative AI for learning and self-development. His work has been published in Nature Machine Intelligence, ACM SIGCHI, and SIGGRAPH, PNAS, and featured in outlets such as The New York Times, Scientific American, and MIT Tech Review. Recognized by TIME's Best Inventions of 2023 and Fast Company's World Changing Ideas, his projects have been exhibited internationally and he has collaborated with NASA, OpenAI, Microsoft Research, and others. He also co-designed one of MIT's first courses on Generative AI and co-created Netflix's 2024 sci-fi anthology Tomorrow and I.

**Anoop Sinha, Ph.D.** is currently a Research Director focusing on AI and Future Technologies at Google, where he leads research into new interfaces and previously directed cross-company AI initiatives involving data and development. Holding a Ph.D. from the University of California, Berkeley, his career spans significant leadership roles in AI and HCI across major tech companies, including Head of Siri ML and Knowledge at Apple and Sr. Applied Research Scientist Manager at Meta (FAIR X - HCI & AI). His expertise lies at the intersection of machine learning, human-computer interaction, search quality, and knowledge representation.

**Yike (Cassandra) Shi** is a Research Associate jointly affiliated with New York University and Carnegie Mellon University. Her research focuses on requirements-driven prompting for LLMs, where she developed a web-based system that compiles user requirements into optimized prompts and integrates automated evaluation mechanisms. She also has industry experience as an AI Infrastructure Intern at DeepLang AI, where she optimized CUDA kernels and improved inference efficiency for quantized models. Shi's academic projects span speech recognition, face recognition, and Transformer-based speech-to-text systems, as well as systems-level programming, database design, and game development. She has been recognized on the Dean's List across multiple semesters.

**Jenny T. Liang** is a PhD student in Software Engineering at Carnegie Mellon University, advised by Brad A. Myers. Her research sits at the intersection of software engineering, HCI, and applied machine learning, focusing on how developers interact with AI-powered tools and how to design more usable systems. She has published in leading venues such as ICSE, FSE, and CHI, receiving awards including the ACM SIGSOFT Distinguished Paper Award. In addition to her research, Jenny has been active in community-building — organizing workshops at ICSE and CHI that bring together researchers across software engineering, HCI, and AI. She also has industry experience through internships at Apple, Microsoft, and AI2, and is dedicated to mentoring and service within the academic community.

**Lama Ahmad, Ph.D.** is a researcher and technology professional currently leading partnerships and research on the risks and social impacts of AI at OpenAI's Safety Systems team. She also

serves as a Term Trustee for The Asia Foundation, guiding governance and strategy. Previously, Lama worked on Facebook's Open Research & Transparency team, focusing especially on democracy, elections, and the societal consequences of social media platforms. During her Luce Scholar year, she studied the ethics of data-driven technologies at the U.N. Global Pulse Lab in Jakarta, applying human-centered design across Southeast Asia. She is a passionate advocate for equity, inclusion, and interdisciplinary approaches in tech and policymaking.

**Tanu Mitra, Ph.D.** is an Associate Professor in the Information School at the University of Washington, with an affiliate appointment in the Paul G. Allen School of Computer Science & Engineering. She is also the Founding Co-Director of RAISE (Responsibility in AI Systems and Experiences). Her research focuses on Human-Centered AI and Responsible AI, combining computational techniques, NLP, and social science principles to study human behavior and interaction in large-scale online systems. An interdisciplinary scholar, Mitra employs methods from HCI, data science, and AI to design systems that foster responsible and effective human-computer and human-human communication. Prior to UW, she was an Assistant Professor in Computer Science at Virginia Tech and received her Ph.D. in Computer Science from Georgia Tech.

**Brad A. Myers, Ph.D.** is the Charles M. Geschke Director of the Human-Computer Interaction Institute and Professor in the School of Computer Science at Carnegie Mellon University, with an affiliated appointment in the Software and Societal Systems Department. He is an ACM Fellow, IEEE Life Fellow, CHI Academy member, and recipient of the 2017 ACM SIGCHI Lifetime Achievement Award in Research. His book, *Pick, Click, Flick!* won the 2025 CBI HCI History Prize. Myers has authored or edited more than 550 publications, including three books, with 19 Best Paper Awards and 6 Most Influential Paper Awards. He has consulted for over 90 companies on UI design and regularly teaches HCI and software design. His research spans interaction techniques, developer experience, API usability, end-user software engineering, programming by example, and visual programming. He has helped organize and run multiple workshops and conferences.

**Yang Li, Ph.D.** is a Senior Staff Research Scientist at Google DeepMind and Affiliate Associate Professor at the University of Washington CSE. His research lies at the intersection of HCI and AI, with a focus on user interface understanding, automation, generation, and code generation for UIs and apps. He has advanced deep learning methods such as Fourier Positional Encoding and area attention, and created influential benchmarks like screen2words and seq2act. Li pioneered on-device interactive ML on Android, leading to features like next app prediction and Gesture Search. He has published widely across HCI and ML venues (CHI, UIST, ICML, NeurIPS, ICLR, ACL, CVPR) and received multiple Best Paper and Lasting Impact Awards. An ACM Distinguished Scientist, Li co-edited AI for HCI: A Modern Approach and organized the inaugural AI & HCI workshops at ICML. He earned his Ph.D. in Computer Science from the Chinese Academy of Sciences and completed postdoctoral research at UC Berkeley EECS.

# 4 Workshop Schedule and Activities

We propose a long, in-person workshop with 180 minutes (including breaks). The workshop is designed to *balance knowledge sharing, interactive discussions, and collaborative activities*, allowing participants to meaningfully connect with others in the AI alignment community. The tentative workshop schedule is detailed in Table 1.

## 4.1 Before the Workshop

Before the workshop session, we will invite participants through social media promotion and professional mailing lists. We expect attendees from diverse backgrounds with varying levels of familiarity and seniority with the topic. To facilitate pre-workshop engagement, we previously created a BiAlign Slack channel (now with 250+ participants), where attendees share materials, papers, and networking opportunities. For CHI 2026, we will expand these platforms to support both pre-event collaboration and post-workshop discussions.

## 4.2 Workshop Schedule

The workshop will be held for 180 min with breaks in the afternoon and is organized into two main sessions. We describe more details about the activities below.

**Keynote Talks:** Two distinguished scholars anchor our program. Dr. Lama Ahmad (OpenAI) will open the first session with a keynote that connects alignment research to human-centered values, setting the tone for the day. Dr. Elizabeth F. Churchill (MBZUAI) will kick off the second session with a forward-looking talk, drawing from her expertise in HCI and AI to inspire cross-disciplinary dialogue. Each keynote will include space for questions, ensuring participants can engage directly with the speakers.

**Paper Presentations:** Accepted papers will be shared through *lightning talks*, giving authors an opportunity to present their ideas in a lively, fast-moving format that sparks curiosity and discussion. A *poster session* later in the workshop will encourage one-on-one exchanges, deeper conversations, and networking across disciplines, providing ample opportunities for participants to find shared interests.

**Group Activity 1 | Concept Mapping & Solution Ideation:** The first collaborative session centers on a *concept mapping and solution ideation* exercise. Working in groups, participants will identify challenges, connect ideas across disciplines, and co-develop creative solutions. This interactive mapping process ensures that every participant's perspective contributes to a shared vision.

**Group Activity 2 | On-the-spot Paper Writing:** Building on the momentum, the second group session introduces an *on-the-spot paper writing* challenge. Teams will synthesize earlier discussions into short outlines or position pieces, capturing fresh insights in a format that can lead to concrete post-workshop collaborations. This activity not only encourages creativity but also creates tangible outputs participants can take forward.

**Insight Sharing & Closing Remarks:** The workshop concludes with *Closing Remarks*, summarizing takeaways and outlining next steps for sustained engagement. Participants are encouraged to continue discussions and collaboration through the dedicated *Slack Channel* and workshop website. This structure is designed to foster

| Time | Activity |
| --- | --- |
| **Session 1 (90min)** | |
| 15 min | Welcome & Overview |
| 20 min | Keynote Talk 1: Lama Ahmad (OpenAI) |
| 20 min | Lightning Talks by Authors |
| 35 min | Group Activity 1: Concept Mapping & Solution Ideation |
| **Session 2 (90min)** | |
| 20 min | Keynote Talk 2: Elizabeth F. Churchill (MBZUAI) |
| 20 min | Poster Session & Networking |
| 30 min | Group Activity 2: On-the-spot Paper Writing |
| 20 min | Insight Sharing & Closing Remarks |

**Table 1: Workshop schedule (180-min long session) with papers, posters, group activities, and discussions.**

active participation, collaboration, and networking, while allowing participants to explore human-AI alignment topics in depth.

### 4.3 Post-Workshop & Plans to Publish Proceedings

We plan to compile a comprehensive report summarizing key discussions, presentations, and findings, which will be shared via open-access platforms such as ArXiv and the workshop website. Outcomes from the On-the-spot Paper Writing session may be developed into full papers for submission to HCI venues such as CHI.

**Plans to Publish Proceedings:** We plan to publish workshop proceedings by collecting accepted papers and curating an edited volume, special journal issue (e.g., ACM ToCHI or ACM TIIS), or as online proceedings via platforms such as CEUR-WS. This will ensure the workshop outcomes reach a broad and sustained audience while supporting the continued development of the AI alignment research community.

**Offline Materials:** Offline and asynchronous access will be provided through the workshop website and Slack, including the program schedule, list of organizers and speakers, pre-prints of accepted papers, and recordings of presentations, which will also be made publicly available on YouTube.

### 4.4 Intended Community & Expected Size

We expect workshop attendees to include academic and industry researchers and practitioners broadly interested in AI alignment topics, coming from diverse disciplines such as human-computer interaction, AI, machine learning, psychology, and social sciences, without requiring deep technical expertise in AI.

**Interdisciplinary Community Connections:** To foster interdisciplinary community connections, we will also host a Machine Learning-oriented Bidirectional Human-AI Alignment Workshop, and our shared Slack Channel enables ongoing collaboration among researchers from HCI and ML communities, and other alignment-focused communities. To ensure meaningful, in-depth discussions, the workshop is tailored for **30-50** in-person participants. We may adjust the structure to accommodate a slightly larger group while preserving interactive engagement.

## 5 Logistics and Accessibility

We are committed to creating an inclusive workshop environment for all participants, including those with cognitive, mental health, or physical disabilities. Authors will be encouraged to make their position papers accessible, and guidance will be provided for accepted papers, including alt-text for images and tables and clear document structure for screen readers. During the workshop, participants will be asked to follow accessibility best practices, such as using microphones and enabling captions for all presentations.

## 6 Call For Participation

We invite researchers and practitioners from academia and industry to join our Bidirectional Human-AI Alignment (BiAlign) Workshop at CHI 2026. As AI systems increasingly permeate everyday life, alignment requires dynamic, reciprocal processes in which humans and AI adapt to each other over time. This workshop provides an interactive forum to explore value-centered alignment, human-AI interaction design, evaluation methods, and strategies for dynamic co-evolution. The workshop will feature paper presentations, poster sessions, and collaborative group activities such as on-the-spot paper writing, concept mapping, and solution ideation. These activities are designed to foster interdisciplinary knowledge creation, critical discussion, and co-development of new research directions. We welcome submissions of position papers, posters, or brief research notes that address human-AI alignment from HCI, AI, psychology, social sciences, or related domains. Accepted participants are expected to attend the workshop, with at least one organizer per accepted submission present. Key workshop topics include:

- **Value-Centered Alignment Objectives:** Embedding fairness, agency, equity, and responsibility into AI systems
- **Designing and Interacting for Alignment:** Interfaces, explanation, and participatory methods for steering AI
- **Evaluating Alignment and Societal Impacts:** Metrics and frameworks for technical and social assessment
- **Dynamic Co-Evolution of Human-AI Futures:** Strategies to maintain alignment as humans-AI mutually adapt

We expect 30-50 in-person participants, ensuring a highly interactive environment. We welcome participants from HCI, CSCW,

AI, design, psychology, communication, and policy. Join us to connect with the growing BiAlign community, engage in hands-on collaborative activities, and help shape reciprocal and responsible human-AI futures.

# References

[1] Micah Carroll, Alan Chan, Henry Ashton, and David Krueger. 2023. Characterizing manipulation from AI systems. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*. 1–13.

[2] Micah Carroll, Davis Foote, Anand Siththaranjan, Stuart Russell, and Anca Dragan. 2024. AI Alignment with Changing and Influenceable Reward Functions. *arXiv:2405.17713* (2024).

[3] 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).

[4] Valdemar Danry, Pat Pataranutaporn, Yaoli Mao, and Pattie Maes. 2023. Don't just tell me, ask me: Ai systems that intelligently frame explanations as questions improve human logical discernment accuracy over causal ai explanations. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–13.

[5] Cathy Mengying Fang, Auren R Liu, Valdemar Danry, Eunhae Lee, Samantha WT Chan, Pat Pataranutaporn, Pattie Maes, Jason Phang, Michael Lampe, Lama Ahmad, et al. 2025. How AI and Human Behaviors Shape Psychosocial Effects of Extended Chatbot Use: A Longitudinal Randomized Controlled Study. *arXiv preprint arXiv:2503.17473* (2025).

[6] Batya Friedman. 1996. Value-sensitive design. *interactions* 3, 6 (1996), 16–23.

[7] Mitchell L Gordon, Michelle S Lam, Joon Sung Park, Kayur Patel, Jeff Hancock, Tatsunori Hashimoto, and Michael S Bernstein. 2022. Jury learning: Integrating dissenting voices into machine learning models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*.

[8] 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.

[9] Alicia Guo, Pat Pataranutaporn, and Pattie Maes. 2024. Exploring the impact of AI value alignment in collaborative ideation: Effects on perception, ownership, and output. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–11.

[10] Maurice Jakesch, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. 2023. Co-writing with opinionated language models affects users' views. In *Proceedings of the 2023 CHI conference on human factors in computing systems*. 1–15.

[11] Yuyang Jiang, Longjie Guo, Yuchen Wu, Aylin Caliskan, Tanushree Mitra, and Hua Shen. 2018. Beyond One-Way Influence: Bidirectional Opinion Dynamics in Multi-Turn Human-LLM Interactions. (2018).

[12] Sayash Kapoor, Rishi Bommasani, Kevin Klyman, Shayne Longpre, Ashwin Ramaswami, Peter Cihon, Aspen Hopkins, Kevin Bankston, Stella Biderman, Miranda Bogen, et al. 2024. On the societal impact of open foundation models. *arXiv preprint arXiv:2403.07918* (2024).

[13] Auren R Liu, Pat Pataranutaporn, and Pattie Maes. 2025. The Heterogeneous Effects of AI Companionship: An Empirical Model of Chatbot Usage and Loneliness and a Typology of User Archetypes. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, Vol. 8. 1585–1597.

[14] Qianou Ma, Hua Shen, Kenneth Koedinger, and Sherry Tongshuang Wu. 2024. How to teach programming in the ai era? using llms as a teachable agent for debugging. In *International Conference on Artificial Intelligence in Education*. Springer, 265–279.

[15] Subhankar Maity and Manob Jyoti Saikia. 2025. Large Language Models in Healthcare and Medical Applications: A Review. *Bioengineering* 12, 6 (2025), 631.

[16] 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.

[17] Pat Pataranutaporn. 2024. *Cyborg Psychology: The Art & Science of Designing Human-AI Systems that Support Human Flourishing*. Massachusetts Institute of Technology.

[18] Pat Pataranutaporn, Ruby Liu, Ed Finn, and Pattie Maes. 2023. Influencing human–AI interaction by priming beliefs about AI can increase perceived trustworthiness, empathy and effectiveness. *Nature Machine Intelligence* 5, 10 (2023), 1076–1086.

[19] Snehal Prabhudesai, Ananya Prashant Kasi, Anmol Mansingh, Anindya Das Antar, Hua Shen, and Nikola Banovic. 2025. " Here the GPT made a choice, and every choice can be biased": How Students Critically Engage with LLMs through End-User Auditing Activity. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–23.

[20] 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.

[21] Jérémy Scheurer, Mikita Balesni, and Marius Hobbhahn. 2023. Large language models can strategically deceive their users when put under pressure. *arXiv preprint arXiv:2311.07590* (2023).

[22] Hua Shen, Nicholas Clark, and Tanushree Mitra. 2025. Mind the Value-Action Gap: Do LLMs Act in Alignment with Their Values? *arXiv preprint arXiv:2501.15463* (2025).

[23] Hua Shen, Tiffany Knearem, 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).

[24] Hua Shen, Tiffany Knearem, Reshmi Ghosh, Michael Xieyang Liu, Andrés Monroy-Hernández, Tongshuang Wu, Diyi Yang, Yun Huang, Tanushree Mitra, Yang Li, et al. 2025. Bidirectional Human-AI Alignment: Emerging Challenges and Opportunities. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–6.

[25] Hua Shen, Tiffany Knearem, 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).

[26] Hua Shen, Ziqiao Ma, Reshmi Ghosh, Tiffany Knearem, Michael Xieyang Liu, Tongshuang Wu, Andrés Monroy-Hernández, Diyi Yang, Antoine Bosselut, Furong Huang, et al. [n. d.]. ICLR 2025 Workshop on Bidirectional Human-AI Alignment. In *ICLR 2025 Workshop Proposals*.

[27] 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).

[28] **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)*.

[29] Antti Väänänen, Keijo Haataja, Katri Vehviläinen-Julkunen, and Pekka Toivanen. 2021. AI in healthcare: A narrative review. *F1000Research* 10 (2021), 6.

[30] 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].

[31] 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.

[32] Hanyi Xu, Wensheng Gan, Zhenlian Qi, Jiayang Wu, and Philip S. Yu. 2024. Large Language Models for Education: A Survey. *arXiv preprint arXiv:2405.13001* (2024). https://arxiv.org/abs/2405.13001