

Ph.D. Statement of Purpose

section: introduction

Despite the impressive capabilities of current AI models, leveraging AI in human-AI teams and human-centered task is yet an open problem, one I am deeply interested in. As a predoctoral researcher at University of Chicago, my research focuses on **human-centered machine learning (HCML)**; in particular, we are designing AI-assistance frameworks for medical training. My research and experience had inspired my interests in the following problems: I am especially interested in (1) designing AI assistance that inspires humans' appropriate reliance instead of blind trust (2) AI learning human perception to provide better AI assistance. I intend to continue working on these problems with the larger goal of building AI assistance that benefit humans while retaining human agency.

Finding my research interest

section: research experience

My research journey began with my work with Prof. Chenhao Tan at UChicago, where we designed a human-compatible model for case-based decision support [1].

subsection: research problem description

In AI explanations, example-based explanations retrieve the nearest neighbors to a test instance, but the similarity metric for retrieval remains understudied. In fact, the metric is often a distance on some AI model embedding, so the retrieved examples may not be informative to humans. Towards this end, we developed a human-compatible model that learns both classification and human perception judgement. Such a model produced representations more aligned with human similarity, thus providing more effective nearest-neighbor explanations, or case-based decision-support.

subsection: specific contributions

My contributions include running synthetic experiments: using an artificial dataset with controllable features, I ran experiments with varying simulated humans, generated by tuning feature weights, and varying decision boundaries, showing our human-compatible model resulted in better decision-support than ML model baselines. With positive results from synthetic experiments, we moved on to real humans. I conducted a human study on a chest X-ray dataset with Prolific crowdworkers, showing our model also provides effective decision support for pneumonia diagnosis.

subsection: research outcomes

The work led to a publication to the Workshop on Human-Machine Collaboration and Teaming in ICML 2022 as well as an under-review submission to a major ML conference.

subsection: new skills

Besides technical skills like designing synthetic experiments and human studies, I learned of an unexpected but useful research practice: be flexible in problem formulation. Our initial goal was to design an AI-driven tutorial for radiology training, largely a machine teaching problem. However, in our experimentations with modeling human learners, we found that a neural network that learned both classification and human similarity judgement produced an interesting representation, one that encoded patterns from human perception. We did not know how to leverage our human-compatible model for teaching, but we thought it may be useful for case-based decision support. Thus we detoured to a different research problem, completed it a project on it, and now we return to the problem of AI-driven tutorial.

section: development of research interest

In this project, we encountered the tradeoff between human agency and performance: showing neutral and informative explanations would retained human agency, but results in lower task performance than showing persuasive explanations that nudged humans to follow the AI. I was intrigued by this dilemma and was drawn to the complexities of human-AI interaction.

subsection: old perspective

Coming from a machine learning background, my perspectives towards research were also largely influenced by ML. Specifically, I treated the human-AI system as another ML model: I thought that if we could build a more generalizable model, produce better explanations and achieve higher performance on some metric, we would achieve better human-AI collaboration.

However, through a Human-Centered Machine Learning course and a Human+AI conference, both at UChicago, I gradually formed a new, more thorough perspective towards HAI. I first learned of the many "non-ML" factors that could affect human-AI team performance. On the human side, AI assistance could alter humans' decision-making process and induce bias [2]; lack of human cognitive engagement could cause overreliance or underreliance towards AI [3]. On the AI side, explanations may not align with humans' decision-making process [4]; the presentation of assistance is also crucial, as showing too little or too much information may be problematic.

More important, I learned of the misalignments between HAI research and real life scenarios, including misalignment in the task and objective [5], in discredita of the assistance, in the subjects that use our assistance.

A more important lesson I learned was a more human-centered perspective: to take a step back and question whether the task at hand is appropriate for a real-life scenario that humans use.

future directions

My experiences collectively shaped my research interests in improving human-AI interaction and collaboration. Besides my interest, I am also very motivated in this direction of research because of its interdisciplinary nature and the visible impact that results from research. I hope to continue this line of work in my PhD. Currently, I am interested in the following problems:

1. How can we design AI assistance that inspires humans' appropriate reliance?

In our project, we spent a large amount of time devising a decision support system that differed from existing literature providing model explanations: our goal was to provide faithful and neutral evidence for humans to make independent decisions. Thus we eventually devised a neutral decision support policy that did not reveal AI's predicted label and provided nearest-neighbor explanations from all classes. We showed such a policy provided effective decision support, but it was less effective than a persuasive decision support policy closer to model explanations.

This alludes to a larger problem: the tradeoff between human agency and human-AI team performance. Many human-AI teams provide humans with "persuasive" AI assistance that improves humans performance but at the expense of humans' overreliance and blind trust towards AI [6]. This is undesired and unethical, especially in high-stake domains where humans should have the last say. On the other hand, many work [6], including ours, also show that neutral, non-persuasive assistance that dissuade humans from blindly following AI perform no better than persuasive assistance. Thus, I am interested in solving this dilemma.

2. AI learning human perception to provide better AI assistance

As our work showed, AI learning human perception can provide better decisions support performance; I believe in general, incorporating more human knowledge and intuition to AI frameworks has potential to provide more effective AI assistance, but many issues remain unsolved in this nascent direction. For example, our work focuses on visual perception on images as it is the easiest form of perception for humans; perception for other modalities like audio and text are also important but much more difficult to collect and utilize. The limitations and bounds of AI learning human perception is also unexplored, including the amount of data needed, the format of data (we used triplet annotations), how well neural networks can learn, etc.

References

- [1] Han Liu, **Yizhou Tian**, Chacha Chen, Shi Feng, Yuxin Chen, and Chenhao Tan. Towards effective case-based decision support with human-compatible representations. *Proceedings of the 1st ICML 2022 Workshop on Human-Machine Collaboration and Teaming*, 2022.
- [2] Ben Green and Yiling Chen. Algorithmic risk assessments can alter human decision-making processes in high-stakes government contexts. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):1–33, October 2021.
- [3] Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z. Gajos. To trust or to think: Cognitive forcing functions can reduce overreliance on ai in ai-assisted decision-making. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–21, April 2021.
- [4] Yaniv Yacoby, Ben Green, Christopher L. Griffin Jr., and Finale Doshi-Velez. “if it didn’t happen, why would i change my decision?”: How judges respond to counterfactual explanations for the public safety assessment. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, 10(1):219–230, Oct. 2022.
- [5] Luke Guerdan, Kenneth Holstein, and Zhiwei Steven Wu. Under-reliance or misalignment? how proxy outcomes limit measurement of appropriate reliance in ai-assisted decision-making. *CHI TRAIT ’22: Workshop on Trust and Reliance in AI-Human Teams*, 2022.
- [6] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, CHI ’21, New York, NY, USA, 2021. Association for Computing Machinery.