# Towards Better Robot Learners: Leveraging Implicit and Explicit Human Feedback Together in Human-Robot Interactions

## Kate Candon

Yale University, New Haven, CT, USA 06511
kate.candon@yale.edu

## Abstract

My work aims to enable robots to better learn from human feedback in human-robot interactions. The way in which people want to collaborate with a robot can vary person-to-person, interaction-to-interaction, or even within an interaction with a given person. Thus, robots need to be able to adapt their behavior during interactions. Robots typically learn from humans via explicit feedback, such as evaluative feedback, preferences, or demonstrations. We know that humans also provide additional information implicitly through non-verbal behavior that gives clues about their internal states during interactions. My work investigates how we can incorporate both kinds of feedback into robot learning paradigms.

## Introduction

Imagine you are making pizza with a robot. The robot already knows generally how to make a pepperoni pizza, but does not know your personal preferences, such as your preference to add pepperoni before cheese. If the robot first passes you tomato sauce after you have rolled out the dough, you might provide positive feedback via the "good job" button. If the robot later hands you cheese before it has handed you pepperoni, this would not align with your preferences, but you might not think handing you the cheese prematurely is enough of an error to warrant providing negative feedback via a "bad job" button. However, you might roll your eyes and put the pepperoni off to the side. I believe it would lead to better interactions, and thus sustained usage of robots, if the robot was able to learn from both the explicit (button press) and implicit (eye roll, putting ingredient aside) human feedback together during the interaction. Thus, my research aims to answer the question: **How can we leverage implicit and explicit human feedback together when a robot is learning from a human in human-robot interactions?**

My work aims to enable robots to more effectively learn how to collaborate with people in a personalized manner. As people interact with robots in a range of tasks across a variety of settings, tasks will become less objective and increasingly driven by personal preferences and proclivities (Bıyık, Talati, and Sadigh 2022). It thus becomes advantageous when robots can adapt to individual preferences (Thomaz, Hoffman, and Cakmak 2016).

Typically, robots learn from humans via explicit feedback (Cui et al. 2021), such as evaluative feedback, preferences, or demonstrations. Recent work has investigated how to combine different types of feedback when a robot is learning from a human teacher (Jeon, Milli, and Dragan 2020; Fitzgerald et al. 2022). However, there are challenges with relying solely on explicitly provided feedback. Providing explicit feedback can place cognitive burden on the human interactant, interrupt the flow of the interaction, and take attention away from the person's own tasks during a collaboration. Additionally, it is well known that there are shortcomings with how humans provide feedback when teaching a robot, such as providing less explicit feedback as an interactions progresses (Li et al. 2013), and that these shortcomings present challenges for machine learning approaches.

One approach to combat the challenges posed by relying solely on explicit feedback is to look toward implicitly provided feedback, such as facial reactions, eye gaze, or human actions. Some work has explored how robots can learn from implicit human feedback (e.g., (Cui et al. 2021)). However, interpreting implicit feedback is challenging as communicative signals are highly individualized and can hold different means across scenarios or cultures (Barrett et al. 2019).

Prior work, both my own and that of other researchers, has investigated explicit and implicit human feedback individually. My current work and proposed future work focuses instead on reasoning about both explicit and implicit human feedback together to enhance robots' capabilities to learn to adapt to individual preferences during human-robot interactions compared to existing approaches. My hope is that considering the two types of feedback in conjunction will help to offset some of the shortcomings, thus improving robots' abilities to more consistently act in accordance with user preferences. Improving robots' ability to learn from humans would enable more seamless human-robot interactions.

## Prior Work

As an initial step, I first explored people's perceptions of an agent's helpfulness. In an exploratory study, we examined which factors influenced the perceived helpfulness of an agent in an interactive video game (Candon et al. 2022). Even though participants were given a clear and objective goal, they had different interpretations of whether or not assistive behaviors from the agent were helpful. Perceived

helpfulness was more strongly correlated with emotional perceptions (e.g., how annoying the human reported the agent to be) than game objectives (e.g., how many points the agent scored for the human). These findings highlighted that it is challenging to know what to optimize for when designing a robot's behavior, so it is important that robots can learn to adapt their behavior during interactions.

Robots typically learn via explicit feedback from humans, but we know that there are limitations with how humans provide feedback. Thus, some of my prior work has investigated how we might mitigate challenges with explicit feedback. In another study, we studied how a robot should remind people to provide explicit feedback during a fast-paced, cooperative interaction (Candon et al. 2023b). We found that robots were able to influence participants to provide feedback more quickly and more frequently by reminding them to provide feedback before a change in robot behavior.

My prior work has also investigated the potential of using implicit feedback that humans naturally provide. We found that even without explicitly directing people to be expressive, we could leverage nonverbal behavior to improve our ability to predict their preferences across different agent behaviors in an interactive game (Candon et al. 2023a). In other work, we analyzed the data collected in two separate human-robot interactions and found that interaction history is an important factor that influences human nonverbal reactions to robots (Candon et al. 2024).

## Current Work

While my prior work has investigated explicit and implicit human feedback individually, my current work explores how to incorporate them together into robot learning paradigms. Though I am hoping to increase the amount of feedback provided by humans to a learning algorithm by considering both explicit and implicit human feedback, the overall amount of feedback is still limited by the fact that we are using real humans. This means that a robot must use the communicative signals that a human provides in an effective manner, making learning practical and low effort for the human. In my current work, **I am investigating how to use explicit and implicit feedback to learn a reward function that models human preferences for an interaction via inverse reinforcement learning**. With such a reward function, the robot can then optimize its behavior to satisfy human preferences.

In advance of the workshop, I hope to have preliminary results for an algorithm learning from explicit and implicit feedback together. The setup collects explicit feedback via evaluative button presses and implicit feedback via human actions in the task. I am working on algorithms to update the robot's belief over the human's reward function from observed feedback. I have a simulated task environment with a simulated human for quick prototyping, and hope to run a user study on our real-world experimental setup.

Additionally, I hope to begin training models to recognize and leverage other human nonverbal behavior during the interaction. In particular, I am interested in analyzing the data we collect during a user study to see if facial reactions or human gestures seem to be more useful in predicting when a participant might provide feedback.

## Future Work

In the twelve to eighteen months following the workshop, I plan to investigate two research directions: **RQ1: How can a robot use implicit feedback to decide when to query a human for explicit feedback?** and **RQ2: How can a robot use implicit feedback to qualify explicit feedback?** Once we build a better understanding of how to incorporate both implicit and explicit feedback together, I want to explore how they can be used to make better sense of each other.

## Acknowledgments

## References

Barrett, L. F.; Adolphs, R.; Marsella, S.; Martinez, A. M.; and Pollak, S. D. 2019. Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements. *Psych. science in the public interest*, 20(1).

Bıyık, E.; Talati, A.; and Sadigh, D. 2022. Aprel: A library for active preference-based reward learning algorithms. In *Proceedings of HRI*. IEEE.

Candon, K.; Chen, J.; Kim, Y.; Hsu, Z.; Tsoi, N.; and Vázquez, M. 2023a. Nonverbal Human Signals Can Help Autonomous Agents Infer Human Preferences for Their Behavior. In *Proceedings of AAMAS*.

Candon, K.; Georgiou, N. C.; Zhou, H.; Richardson, S.; Zhang, Q.; Scassellati, B.; and Vázquez, M. 2024. RE-ACT: Two Datasets for Analyzing Both Human Reactions and Evaluative Feedback to Robots Over Time. In *Proceedings of HRI*.

Candon, K.; Hsu, Z.; Kim, Y.; Chen, J.; Tsoi, N.; and Vázquez, M. 2022. Perceptions of the Helpfulness of Unexpected Agent Assistance. In *Proceedings of HAI*.

Candon, K.; Zhou, H.; Gillet, S.; and Vázquez, M. 2023b. Verbally Soliciting Human Feedback in Continuous Human-Robot Collaboration: Effects of the Framing and Timing of Reminders. In *Proceedings of HRI*.

Cui, Y.; Koppol, P.; Admoni, H.; Niekum, S.; Simmons, R.; Steinfeld, A.; and Fitzgerald, T. 2021. Understanding the Relationship between Interactions and Outcomes in Human-in-the-Loop Machine Learning. In *Proceedings of IJCAI*.

Fitzgerald, T.; Koppol, P.; Callaghan, P.; Wong, R. Q. J. H.; Simmons, R.; Kroemer, O.; and Admoni, H. 2022. IN-QUIRE: INteractive querying for user-aware informative REasoning. In *6th Annual Conference on Robot Learning*.

Jeon, H. J.; Milli, S.; and Dragan, A. 2020. Reward-rational (implicit) choice: A unifying formalism for reward learning. In *NEURIPS*, volume 33.

Li, G.; Hung, H.; Whiteson, S.; and Knox, W. B. 2013. Using Informative Behavior to Increase Engagement in the Tamer Framework. In *Proceedings of AAMAS*.

Thomaz, A.; Hoffman, G.; and Cakmak, M. 2016. Computational Human-Robot Interaction. *Foundations and Trends® in Robotics*, 4(2-3).