

Helping Users Develop Accurate Mental Models of Robots’ Perceptual Capabilities: A First Approach

Matthew Rueben¹, Matthew Tang², Eitan Rothberg³, and Maja J. Matarić¹

Abstract—Understanding a robot’s perceptual capabilities is important for users to make informed decisions about their privacy. Privacy, in turn, is tied to trust and is also important for user acceptance of robots. We are beginning a research program that is exploring (1) how a user forms a mental model about a robot’s perceptual capabilities, (2) how the robot can estimate such models via observation, and (3) how the robot can take actions to correct an inaccurate model, i.e., teach the user about its capabilities. We introduce the research topic, discuss recent relevant work, and then describe our initial approach toward enabling robots to model and influence user beliefs.

I. INTRODUCTION

Before interacting with a robot, users ought to know its privacy-relevant capabilities: can it record my voice? Recognize me? Can it see behind itself? The focus of this work is on the robot’s perceptual capabilities. Unfortunately, users can form wrong impressions of robots’ perceptual capabilities due to stereotyped expectations [1], difficult-to-see sensors [2], or potentially misleading robot behaviors [3]. A user’s trust in the robot can be decreased if his or her expectations are violated by the robot’s actions.

In this work we aim to (1) give robots an artificial understanding of users’ beliefs about their (the robots’) perceptual capabilities, (2) enable robots to automatically select actions that will encourage users to adopt more accurate beliefs, and (3) combine those two components in a framework by which the robot decides when to take these actions and reasons about whether they have successfully influenced user beliefs.

II. BACKGROUND: MENTAL MODELS OF ROBOTS

Users cannot see what is going on “inside” a robot. Instead, they rely on the “system image” [4]: the parts of the system that are observable, including how a robot looks and behaves. These observations are used to develop a mental model of the robot that explains and predicts its behaviors [4]. Understanding how users form mental models of robots—especially of their perceptual capabilities [5]—is important for enabling robots that respect user privacy.

III. RELATED WORK

Kwon et al. [6] demonstrated a way to automatically generate actions that express the robot’s (in)capabilities to users who are watching. Inspired by that work, we additionally aim

for the robot to be able to maintain a real-time estimate of the user’s mental model, and use that estimate both to make decisions about *when* to perform informative actions and to check whether those actions are successfully changing user beliefs.

Cha et al. [7] studied how a robot’s behaviors can influence user mental models of its manipulation and social capabilities, but not its *perceptual* capabilities. Rueben et al. explored models of perceptual capabilities in particular and recorded examples of users inferring the robot’s sensing capabilities from its behavior alone [3]. They also observed people mistakenly believing that a robot was fully autonomous when in fact the robot was teleoperated [3]—a mistake with implications for privacy and trust.

This project is also concerned with understanding how users learn about a robot’s perceptual capabilities, but we also aim to build a computational model of this learning process by which an autonomous robot can intelligently influence the user’s beliefs. To our knowledge, this fully autonomous control loop has only been implemented for mental models of robot capabilities by Nikolaidis et al. [8] in the context of robot *manipulation* capabilities. This work focuses on enabling robots to act based on users’ mental models of their *perceptual* capabilities.

IV. APPROACH

Our goal is to develop a way for a robot to automatically model and influence a user’s mental model of its perceptual capabilities. Our initial approach has been to model this problem as a partially observable Markov decision process (POMDP). A POMDP consists of sets of possible states, actions, and observations; transition probabilities between states; probabilities of making different observations; and a reward for each state-action pair. These we define as follows.

In our model, the **state** is the user’s actual belief about the robot’s perceptual capabilities. The robot will maintain a belief distribution of the probability of each belief the user could be holding. The initial belief distribution could be uniform, or it could be informed by what we know about typical first impressions of this robot based on its appearance and the user’s demographics or personality.

The robot’s **action space** depends on the specific application. The **transition function** will encode our understanding of how the user will adjust their belief according to the robot’s behavior. We might obtain this mapping between robot actions and user belief updates experimentally, and eventually explore having the robot learn the mapping from observation.

¹M. Rueben (corresponding author) and M. J. Matarić are with the Department of Computer Science, University of Southern California, Los Angeles, CA 90007. Emails: {mrueben, mataric}@usc.edu

²M. Tang is with the University of California, Berkeley; Berkeley, CA 94720. Email: matthewtang@berkeley.edu

³E. Rothberg is with the Ohio State University, Columbus, OH, 43210. Email: rothberg.10@buckeyemail.osu.edu

The **set of observations** is the user’s action space. More difficult is the **observation function**, which encodes our understanding of what different user actions indicate about that user’s beliefs. The aspects (features) of user behavior to consider, as well as the mapping between these aspects and user beliefs, could be experimentally determined at first and then eventually learned by the robot. The mapping could vary between individuals based on how quickly they learn and individual differences in how they behave.

The **reward** will be the robot’s existing task reward structure, augmented with a second reward if necessary to further encourage the robot to teach the user about its perceptual capabilities. This explicit reward for belief accuracy would be proportional to the distance between the belief that the robot thinks the user has and the robot’s actual abilities. We believe that for some types of tasks the robot’s normal task reward will implicitly encourage the robot to take teaching actions—e.g., in collaborative tasks wherein the user can improve the team’s reward if he or she knows what capabilities the robot has. In such cases, additional rewards for teaching would not be necessary to achieve the belief alignment that is desired for users to be sufficiently well-informed when making decisions about their privacy. Another case wherein we would not use an additional reward for teaching would be if the information being taught were not important beyond its potential for improving task performance—here, knowing the robot’s perceptual capabilities is important for making well-informed decisions about user privacy, so we want to incentivize the robot to teach them if the task does not itself provide enough incentive.

First Steps. In our early, exploratory work, we will use well-structured, collaborative games wherein we can guess a reasonably accurate transition function and observation function, and the game reward structure encourages the robot to learn the user’s beliefs and influence them to be more accurate. Next, we will explore a more realistic scenario in which the user and robot both take continuous actions and it is not so clear what these actions imply about the robot’s perceptual capabilities and about the user’s beliefs.

Searching for the Optimal Robot Policy. Using an exact solver on this POMDP would yield a policy that optimally trades off between gathering information about user beliefs, influencing them through the robot’s actions, and taking actions that more directly earn rewards for the robot. However, solving a POMDP becomes quickly intractable for larger problems, so we are also exploring approximate solvers like SARSOP [9]. We have implemented a modified Monte Carlo Tree Search (MCTS) to find good policies in real-time interactions with users but it is still being tested.

V. BROADER VISION

Our goal is to enable robots to intelligently include information about their perceptual capabilities into their otherwise task-focused actions so human coworkers and bystanders can be informed when making decisions about their privacy. More broadly, this work could be extended to other collaborative HRI domains in which the robot has

knowledge that the user lacks. This could contribute to the broader goal of facilitating the development of shared mental models by human-robot teams [10]. Mental model alignment in general could increase users’ trust in the robot, and robot behavior in accordance with human values such as privacy could increase robot acceptance.

ACKNOWLEDGMENTS

This work was supported by an NSF NRI grant for “Socially Aware, Expressive, and Personalized Mobile Remote Presence: Co-Robots as Gateways to Access to K-12 In-School Education”, NSF IIS-1528121. The second and third authors were also partially supported by the USC Viterbi Summer Undergraduate Research Experience (SURE) Program. The authors are grateful to Prof. Stefanos Nikolaidis for his advice on problem formulation. We also acknowledge Christopher Birmingham for his approach to how robots can estimate and influence user trust, which inspired this work in the context of privacy.

REFERENCES

- [1] N. M. Richards and W. D. Smart, “How Should the Law Think About Robots?” in *Proceedings of “We Robot”: The Inaugural Conference on Legal and Policy Issues Relating to Robotics*, Coral Gables, FL, 2012.
- [2] M. K. Lee, K. P. Tang, J. Forlizzi, and S. Kiesler, “Understanding users’ perception of privacy in human-robot interaction,” in *6th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 2011, pp. 181–182.
- [3] M. Rueben, J. Klow, M. Duer, E. Zimmerman, J. Piacentini, M. Browning, F. J. Bernieri, C. M. Grimm, and W. D. Smart, “Mental models of a mobile shoe rack: Exploratory findings from a long-term, in-the-wild study,” *ACM Transactions on Human-Robot Interaction (THRI)*, in Review. Submitted Feb 5, 2019.
- [4] D. Norman, *The design of everyday things: Revised and expanded edition*. Constellation, 2013.
- [5] B. Schafer and L. Edwards, “I spy, with my little sensor: fair data handling practices for robots between privacy, copyright and security,” *Connection Science*, vol. 29, no. 3, pp. 200–209, 2017.
- [6] M. Kwon, S. H. Huang, and A. D. Dragan, “Expressing robot incapability,” in *Proceedings of the 2018 ACM/IEEE International Conference on Human-Robot Interaction*. ACM, 2018, pp. 87–95.
- [7] E. Cha, A. D. Dragan, and S. S. Srinivasa, “Perceived robot capability,” in *24th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*. IEEE, 2015, pp. 541–548.
- [8] S. Nikolaidis, S. Nath, A. D. Procaccia, and S. Srinivasa, “Game-theoretic modeling of human adaptation in human-robot collaboration,” in *12th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 2017, pp. 323–331.
- [9] H. Kurniawati, D. Hsu, and W. S. Lee, “Sarsop: Efficient point-based pomdp planning by approximating optimally reachable belief spaces,” in *Robotics: Science and systems*, vol. 2008. Zurich, Switzerland, 2008.
- [10] M. Scheutz, S. A. DeLoach, and J. A. Adams, “A framework for developing and using shared mental models in human-agent teams,” *Journal of Cognitive Engineering and Decision Making*, vol. 11, no. 3, pp. 203–224, 2017.