

# Social Feedback For Robotic Collaboration

Author Names Omitted for Anonymous Review. Paper-ID [add your ID here]

**Abstract**—Robotic collaboration requires not just communication from the human to the robot, but also from the robot to the human. We call this robot-to-human communication *social feedback*. In order to flexibly and intelligently generate social feedback, we describe a partially Observable Markov Decision Process that incorporates the human’s belief about the robot’s intent. Doing so incentivizes the robot to communicate its current state in order to allow the human participant to make informed decisions about their next actions. Initial evaluation in simulation shows improvement in interaction length in a toy domain while maintaining near-perfect accuracy.

## I. INTRODUCTION

Collaboration is a process that relies heavily on communication for its success. When humans collaborate, this communication is both obvious and implicit—we not only instruct and request aid from each other, but we also tacitly monitor our partners for signs of approval and understanding while producing these signals ourselves. However, in robotics, many of these vital components of successful communication are missing or lost. The absence of these signals likely account for numerous failures in human-robot collaborative tasks. In order to give robots these missing communicative skills, we employ *social feedback* signals to provide human-like communication to both inform the human participant of the robot’s state and request information from the participant.

Feedback refers to the responses that are received by an agent when it takes some action. These feedback responses inform the agent about the success or failure of its actions. *Social feedback* therefore refers to social signals that convey this information. These can be explicit signals, such as “I want this object,” or implicit signals, such as a perplexed expression. When humans interact with another human, they use this feedback to improve the flow and clarity of the collaboration, which in turn improves its success. In this paper, we describe a framework that describes how robotic agent should likewise generate social feedback for its human partner in human-collaborative tasks, and show that it improves the speed and accuracy of human-robot interaction.

Many existing works focus on the task of interpreting human communications[? ? ? ?], but provide no response if the interpretation fails. As a consequence, it becomes difficult to tell whether the robot is still making progress in interpreting the request, or has failed. Some works do allow for robotic responses to ambiguous situations, but the manner of their responses are hard coded, often limited to generic requests, such as “please repeat the question.” In this work, we provide a means to generate these requests in the form of language and gesture that can flexibly represent the robot’s state as well as inform the human in such a way so that allows the human to help the robot as much as possible.



Fig. 1. The experimental setup from the participant’s perspective: a Baxter robot has six objects within reach. The participant stands in front of the table to provide speech and gesture requests

We will solve this problem by formulating it as a Partially Observable Markov Decision Process (POMDP)[? ], which will allow us to dynamically and flexibly determine how to choose social feedback actions. First, we will describe a two agent model and use this to motivate the construction of a POMDP. The crucial improvement we make over existing robot interaction models is that we maintain a state variable that tracks the human’s interpretations of the actions the robot takes. We will next discuss how we solved this POMDP and other measures we took to allow our system to respond dynamically and fluidly.

**stefie10: swap with following paragraph** Specifically, we address the object delivery task. In this task, a set of objects are laid out on a table within a Baxter robot’s reach. A human participant requests an object from the robot using speech and gesture (pointing). The robot must interpret the human’s speech and gesture and deliver the requested object to the human. The robot then repeats with the remaining objects on the table. This task is achievable without social feedback; the robot need only wait until enough information is given and then deliver the correct object. However, we will show that by adding social feedback actions, such as asking questions, looking at objects, and pointing at objects, we will achieve better accuracy and speed as well as improved user experience.

As an initial evaluation, we run our model in simulation to showcase improved speed and accuracy in the object delivery task.

## II. RELATED WORK

**stefie10: search and replace all cite with citet or citep**  
**EXW: The related works are a little out of date; they focus more on multimodal interpretation, even though**

this paper is more about modeling and action generation now

This work is primarily built off of [?] **stefie10: cite icra paper instead**, which describes a system for incremental speech and gesture recognition in an object-delivery domain. In previous work [?], the model from [?] is expanded from a bayesian inference model into a POMDP model that can ask yes-no questions. In this work, we generalize the idea of the affect of robot’s actions on humans state within a POMDP framework.

Work demonstrating the importance of social feedback in human-human communication has been done in the field of psycholinguistics. In [?], one human (labeled the builder) builds a Lego model according to instructions given by another human (labeled the director). In the feedback-free trials, the director’s instructions were prerecorded, and the resulting models were very inaccurate (in fact no model was completely correct). In the feedback trials, errors were reduced by a factor of eight. Our goal is to enable a robot to collaborate with a human in this way.

Other work with collaborative robots exists, for example, [?] have done research with a bar-tending robot. This robot follows a rule-based state estimator, and delivers drinks from fixed positions behind the bar to multiple users based on their speech and torso position. We expand the scope of the problem: we do not use a rule-based state planner, our items are not in fixed positions, and our gesture model uses pointing instead of torso position.

In [?], a robotic building guide directs guests to find specific rooms. Our project addresses a similar domain, requiring the interpretation of users’ requests, but differs in the task and the type of communication necessary to accomplish that task.

Other work involving robotic object delivery also exists. Some approaches have no social feedback and will either deliver the wrong item or do nothing if given a request it does not understand [? ? ? ?]. Language only feedback models also exist [? ? ? ? ? ?], and several gesture only models [? ?].

[?] shows promising work in fusing language and complex gesture to understand references to multiple objects at once. We build off this work by including social feedback.

In the field of computational linguistics, previous work exists in resolving referring expressions incrementally, such as [? ? ?]. Other work in that community also incorporates gesture, and/or eye gaze [? ?], but the given work does not incrementally update gesture along with speech. [?] provides work towards resolving referring expressions in a different domain, but does not address the task of acting on the results of these referring expressions. In [?], they propose a system for planning to ask for clarifications, which covers a wide scope of knowledge failures. In this work, we are interested only in a small subset of these clarifications, and address the problem of how and when these clarifications should be used in a concrete human-robot collaboration task.

POMDP approaches to dialog [?] are quite common, but treat dialog as a discrete, turn-taking interaction. The Dialog State Tracking Challenge [?] a notable driving force for computer dialog understanding, treats dialog in this turn-based way. Our model is also based off a turn-based interaction, but is designed such that it would be suitable for a more incremental approach given the appropriate framework.

Interactive POMDPs (I-POMDPs) [?] describe a multi-agent POMDP that allows agents to anticipate each other’s actions by forming beliefs over action and observation histories. This work attempts to accomplish the same, but makes more specialized assumptions that allow for a more tractable state space.

Alternative approaches to POMDPs include cognitive architectures such as SOAR [?] or DIARC [?]. By taking a probabilistic approach, we can seamlessly fuse information from multiple sources and explicitly reason about the robot’s uncertainty when choosing actions.

### III. TECHNICAL APPROACH

**stefie10: Put a few sentences here giving an overview. no section/subsection without intro text**

#### A. POMDP Overview

Markov Decision Processes (MDPs) are models that describe how an agent can take actions to transition between states, receiving rewards for its actions. Partially Observable Markov Decision Processes (POMDPs) are used to model MDPs where the true state is not known. Instead, the agent receives observations that are generated by the true state, and must infer what the true state is from the observations. Thus, the agent maintains a belief over true states which is updated as it receives new observations. In subsequent sections, we will refer to this belief over hidden states as  $b$ . The agent must then use this belief state to determine which action to take to maximize its expected value of the reward over time. This splits the POMDP into two main components, a state estimator and a policy generator. A typical graphical model for a POMDP is provided in figure ?? **stefie10: cut this figure.**<sup>1</sup>

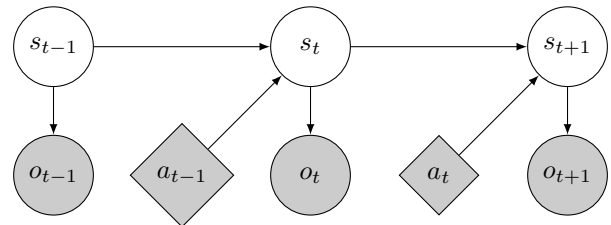


Fig. 2. A graphical model representation of a POMDP

<sup>1</sup>In this model, we have omitted the dependency between an action and the observation from the next timestep. In the POMDP discussed in this paper, we will not use this dependency for any action. Rather, the affect of the actions on the observation is encapsulated in the state.

### B. Model Description

To model this human-robot interactive task, we will use a POMDP. To motivate the construction of our model, consider two-agent model constructed depicted in figure ??.

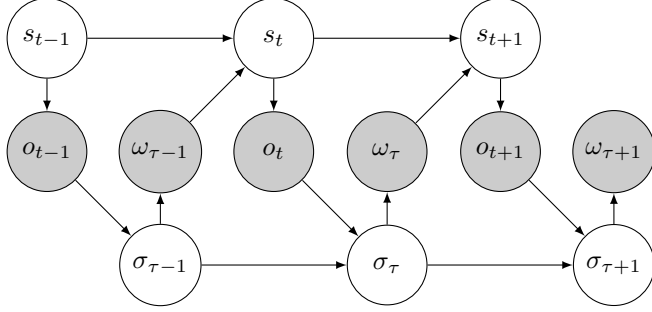


Fig. 3. A two-agent model based off two POMDPs **stefie10: Let's name this model and call it our model.**

We represent human states  $s_t$ , robot states  $\sigma_\tau$ , and observations generated by the robot  $\omega_\tau$  and by the human  $o_t$ . Observe that the structure of this model resembles two POMDPs combined together at their actions and observations. The lower POMDP is the POMDP from the human's perspective: the robot has some hidden states  $\sigma_\tau$ , which the human observes by means of  $\omega_t$ . The human takes actions  $o_t$  to influence the robot's state  $\sigma_t$ . The upper POMDP models the interaction from the robot's perspective: the human's state  $s_t$  is hidden from the robot, and the robot must infer it from observations  $o_t$ . When the robot takes action  $\omega_{\tau-1}$ , it affects the human's state  $s_t$ . Crucially, each agent treats the other agent's action as an observation that influences their belief about the other agent's hidden state. Thus, the human's actions affect the robot's belief about the human's state, which is what we call  $b$  above. Importantly, the reverse is also true: the robot's actions affect the human's belief about what the robot's hidden state is. We will call the human's belief over the robot's hidden state  $\beta$ .

In the following section we will use this dual structure to inform the construction of our POMDP as applied to our object delivery domain.

### C. POMDP Definition

We define our object-delivery POMDP as a tuple  $\{S, A, T, R, \Omega, O\}$ :

- Each  $s \in S$  is a tuple of  $\langle \iota, \beta, \mathcal{I} \rangle$ 
  - $\mathcal{I}$  is the set of all objects that the robot can deliver. Each object is parameterized by a name, a unigram vocabulary, and a position; for example: a red bowl would be represented  $\langle \text{redBowl}, [\text{red}, \text{red}, \text{bowl}, \text{bowl}, \text{plastic}], (1.0, 2.0, 0.0) \rangle$ . We assume the set of all objects is known.
  - $\iota \in \mathcal{I}$  is the object the human desires. This is a hidden variable.
  - $\beta$  is a distribution over the robot's hidden states, as defined above. In this domain, the robot's hidden

state contains which object the robot will hand the human, or which object the robot believes the human wants. We assume the human estimates the distribution over the hidden state by observing the actions the robot has taken. Specifically,  $\beta_t(\iota_t) = p(\iota_t | a_{1:t-1})$ . While this variable is technically unknown, we make some assumptions about the transition functions and initialization so  $\beta_t$ 's value is known at every time step.

Since our state is composed of both hidden and known states, it resembles a Mixed Observability MDP [?], and we leverage the computational benefits of having both hidden and known states.

- The set of actions  $A$  consists of both social and non-social actions. Non-social actions are `pick(i)` (picking up and delivering an object  $i$ ) and `wait`. Social actions are `point(x)`, pointing at a location  $x$ ; `look(x)`, looking at a location  $x$ ; <sup>2</sup> and `say(p)`, informing the human that the robot believes the desired object has property  $p$ , where  $p$  is an element of some object's vocabulary.
- $T(s, a, s') = p(s' | s, a)$ . We make the assumption that the human's desired object  $\iota$  does not change unless the robot picks up object  $\iota$ . The set of all objects  $\mathcal{I}$  transitions to  $\mathcal{I} \setminus \{i\}$  when the robot chooses action `pick(i)`. For each action in  $A$ , we define a "reverse observation function" that describes our assumptions about how the human imagines the robot generates actions (which the human sees as observations) given a hidden robot state. This allows us to define a transition function for  $\beta$ . The transition dynamics are described in greater detail in section ??.
- $R(s, a, s')$  returns a numeric reward from transitioning from  $s$  to  $s'$  by taking action  $a$ . In this domain, we incentivize our robot to pick the correct object by giving it a +10 reward if it delivers the correct object and a -50 reward for picking the incorrect object. We also give negative rewards for taking various actions: `wait` receives a -1 reward (to incentivize the robot to finish the task quickly); `look(x)` receives a reward of -2; `say(p)` receives a reward of -3; `point(x)` receives a reward of -4. These additional penalties for social actions reflect the penalty for "bothering the user", as well as the time it takes to execute these actions. These reward values were chosen experimentally to result in a high rate of pick accuracy.
- Observation  $o \in \Omega$  represents an observation generated by the human. These are tuples of language and gesture:  $\langle l, g \rangle$ . Language utterance  $l$  is represented by a string of any number of words, obtained by transcribing microphone input using webkit's speech recognition API through Google Chrome. A gesture  $g$  is represented by a vector from the participant's shoulder to their

<sup>2</sup> $x$  is chosen from the set of points that describe the surface of the table. For computational reasons, in this work, we restrict it to the set of object locations

wrist, and all gestures are interpreted as a straight-armed point. This vector is obtained using the Kinect’s tracking software.

- $O(o, s, a) = p(o|s, a)$  describes the probability of seeing an observation  $o$  from the human given their state  $s$  and the robot’s last action  $a$  (though assume all observations are independent of  $a$ ). We choose an observation function that reflects that the human is an agent attempting to communicate which object they desire to the robot, and thus chooses to generate observations that are more likely to result in the robot delivering the correct object. This is gone into more detail in the next section, section ??.

For quick reference, a table of variables is provided in figure ??.

1) *Observation Function:* According to our double-agent model, the human emits observations as though it were an agent interacting with our robotic agent. Thus, we choose an observation model that depends on the human’s belief about the robot’s state,  $\beta$ . Specifically, the human will choose an action according to its estimate that the robot will hand them their desired object. In order to define this function, we will first have to define a base-level observation function.

a) *Base-Level Observation Function:* The base level observation function describes the probability of an observation conditioned only on the object:  $p(o|\iota)$ . For our object delivery domain, we will define two base-level observations, one for language and one for gesture. These observation functions are the ones defined in ? ].

*Speech Model:* Language is interpreted according to a smoothed unigram speech model. An utterance  $l$  is broken down into individual words,  $w \in l$ :

$$p(l|\iota) = \prod_{w \in l} p(w|\iota) = \prod_{w \in l} \frac{\text{count}(w, \iota.\text{vocab}) + \alpha}{|\iota.\text{vocab}| + \alpha|\text{words}|}$$

where  $\text{count}(w, \iota.\text{vocab})$  is the number of times word  $w$  appears in  $\iota$ ’s vocabulary. This allows for repeated words to have greater probability of being spoken.<sup>3</sup>  $\alpha$  is a smoothing value chosen to be 0.05.  $|\text{words}|$  is the total number of words in all object vocabularies.

*Gesture Model:* All gestures are interpreted as a straight armed point. These pointing gestures are selected from a normal distribution centered at the object’s location.

Define the angle between the vector defined by the pointing gesture and the vector from the human’s shoulder to the object  $\iota$  to be  $\theta_\iota$ . The probability of a particular gesture is then

$$p(g|\iota) = \mathcal{N}(\theta_\iota|0, v)$$

where  $v$  is a hand-tuned variance. We choose  $v = 0.4$ .

<sup>3</sup>In the future, object vocabularies will be collected from data; repeated words are therefore meaningful as ways that are commonly used to describe objects.

b) *Posterior Observation Function:* We will use the base-level observation function defined above to define a posterior observation function that considers how the base-level observation function will affect the robot’s belief. Specifically, because the human is a rational agent, they will choose to generate an observation that has the greatest probability of the robot handing them their desired object  $\iota$ .

$$p(o|\iota, \beta) = \begin{cases} 1 & \text{if } o = \text{argmax}_o p(\iota|o) \\ 0 & \text{otherwise} \end{cases}$$

To calculate  $p(\iota|o)$ , we can use Bayes rule:

$$p(\iota|o) = \eta \frac{p(o|\iota)p(\iota)}{\sum_{\iota'} p(o|\iota')p(\iota')}$$

where  $\eta$  is a normalization factor.

Next, we will set  $p(\iota) = \beta(\iota)$ .  $p(\iota)$  describes the robot’s belief in  $\iota$ . This value is represented in the robot’s belief-state vector,  $b$ . But since the human does not know  $b$ , it uses its estimate of  $b$ ,  $\beta$ . Our new expression is

$$p(i|o) = \eta \frac{p(o|\iota)\beta(\iota)}{\sum_{\iota'} p(o|\iota')\beta(\iota')}$$

To account for suboptimal behavior and avoid irregularities in our observation function, we soften our assumptions of optimality such that  $1 - \epsilon$  fraction of the time, the human chooses the optimal action. With probability  $\epsilon$  they pick according to the base-level observation function described above.

$$p(o|\iota, \beta) = \begin{cases} (1 - \epsilon) + \epsilon p(o|\iota) & \text{if } o = \text{argmax}_o p(\iota|o) \\ \epsilon p(o|\iota) & \text{otherwise} \end{cases}$$

Using  $\beta$  as a prior for our observation function incentivizes the robot to take actions that modify  $\beta$  in such a way that increase the probability of useful observations, while in turn decreasing the probability of getting observations that tell the robot things it already knows. This result falls out of our assumption that the human is more likely to provide observations that communicate its intentions to the robot well. Consider the following example:

First, we will define a toy domain. In this toy domain, the states are the following:  $\{AA, AB, BA, BB\}$  as well as  $\beta$ , which is a distribution over these states. The observations are  $\{A\_ , \_A, B\_ , \_B\}$ , which mean “the first character is an A”, “the second character is an A”, “the first character is B”, “and the second character is a B” respectively. These observations are provided by the human. The agent can take actions  $CX$  which is informing the human that the agent believes the Xth character is a C, as well as “picking” the object or waiting. For the base level observation, the human always gives truthful observations, and has equal probability of generating an observation pertaining to a particular character. We will see how this affects the POMDP observation function, which incorporates  $\beta$ .

Consider the following situation:



Variable	Explanation
$s = \langle \iota, \beta, \mathcal{I} \rangle$	A single state, defined by a tuple
$\mathcal{I}$	The set of all objects
$\iota \in \mathcal{I}$	The object desired by the human
$\beta$	A distribution over $\mathcal{I}$ . $\beta_t(i_t) = p(i_t a_{1:t-1})$
$i = \langle \text{name, vocab, location} \rangle$	An object, defined by a name, vocabulary, and location
$a$	A robot action, <code>pick(i)</code> , <code>look(x)</code> , <code>point(x)</code> , or <code>say(x)</code> , observed by human
$o = \langle l, g \rangle$	An observation received by the robot, generated by the human.
$l$	A string of language.
$g$	A pointing gesture.
$\epsilon$	The probability of picking a non-optimal observation.
$L$	An indicator variable that is 1 if language is observed.
$G$	An indicator variable that is 1 if gesture is observed.
$c_L$	A constant that describes the probability of observing language.
$c_G$	A constant that describes the probability of observing gesture.
$\alpha$	A smoothing parameter for the unigram model.
$p(a_t i_{t+1})$	A reverse observation function
$A_t$	A diagonal matrix describing the reverse observation function $p(a_t i_{t+1})$ for all values of $i_{t+1}$
$T$	A matrix describing all transition probabilities $p(\iota_{t+1} \iota_t)$

Fig. 4. Table of Variables

Initially both  $b$ , the robot's belief about the human's desired object and  $\beta$ , the human's belief about the robot's belief, are uniform. The true state is  $AA$ .

$t$	$b$	$\beta$	$a$	$o$
0	[0.25, 0.25, 0.25, 0.25]	[0.25, 0.25, 0.25, 0.25]		

The robot then receives an observation  $A_-$ . The new beliefs are:

$t$	$b$	$\beta$	$a$	$o$
0	[0.25, 0.25, 0.25, 0.25]	[0.25, 0.25, 0.25, 0.25]		
1	[0.50, 0.50, 0, 0]	[0.25, 0.25, 0.25, 0.25]		$A_-$

Next, the robot can choose to take an action. If it chooses to wait, this is the resulting state

$t$	$b$	$\beta$	$a$	$o$
0	[0.25, 0.25, 0.25, 0.25]	[0.25, 0.25, 0.25, 0.25]		
1	[0.50, 0.50, 0, 0]	[0.25, 0.25, 0.25, 0.25]		$A_-$
2	[0.50, 0.50, 0, 0]	[0.25, 0.25, 0.25, 0.25]	wait	

Examine the probabilities of each observation.

$$\begin{aligned}
p(A_-|AA, \beta) &= \frac{p(A_-)\beta(AA)}{\sum_{h \in \{AA, AB, BA, BB\}} p(A_-|h)} \\
&= \frac{0.5 * 0.25}{0.5 * 0.25 + 0.5 * 0.25 + 0 + 0} \\
&= 0.5
\end{aligned}$$

$$\begin{aligned}
p(_A|AA, \beta) &= \frac{p(_A)\beta(AA)}{\sum_{h \in \{AA, AB, BA, BB\}} p(_A|h)} \\
&= \frac{0.5 * 0.25}{0.5 * 0.25 + 0 + 0.5 * 0.25 + 0} \\
&= 0.5
\end{aligned}$$

The probabilities of all other actions are 0, since we only give truthful observations.

Notice that this situation is not ideal. According to the human, producing each of these observations equally optimal, so we are equally likely to see either observation in this situation, even though the agent already knows that the first character is an A. Now, consider what would happen if the robot chose the action  $A0$ . We would get the following belief states:

$t$	$b$	$\beta$	$a$	$o$
0	[0.25, 0.25, 0.25, 0.25]	[0.25, 0.25, 0.25, 0.25]		
1	[0.50, 0.50, 0, 0]	[0.25, 0.25, 0.25, 0.25]		$A_-$
2	[0.50, 0.50, 0, 0]	[0.5, 0.5, 0, 0]	$A0$	

If we examine the probabilities again:

$$\begin{aligned}
p(A_-|AA, \beta) &= \frac{p(A_-)\beta(AA)}{\sum_{h \in \{AA, AB, BA, BB\}} p(A_-|h)} \\
&= \frac{0.5 * 0.5}{0.5 * 0.5 + 0.5 * 0.5 + 0 + 0} \\
&= 0.5
\end{aligned}$$

$$\begin{aligned}
p(_A|AA, \beta) &\propto \frac{p(_A)\beta(AA)}{\sum_{h \in \{AA, AB, BA, BB\}} p(_A|h)} \\
&= \frac{0.5 * 0.5}{0.5 * 0.5 + 0 + 0 + 0} \\
&= 1
\end{aligned}$$

Now, observation  $_A$  is the most optimal, so there is high probability the human will generate this observation

for the robot. This observation will give the robot the information it needs to pick the correct object,  $AA$ . This provides an incentive for the robot to choose an action that conveys something about its current state of belief over doing nothing.

*c) Modification for Object Delivery Domain:* In the object delivery domain, we have both speech and gesture, which we assume are conditionally independent given the state.

$$p(o|s) = p(l, g|s) = p(l|s)p(g|s)$$

In addition, it is possible to observe no speech or no gesture input. Let  $L$  be a random variable that is 1 if the agent receives a language observation and 0 otherwise. Similarly, let  $G$  be a random variable that is 1 if the agent receives a gesture observation and 1 otherwise. For all states  $s$ :

$$\begin{aligned} p(L = 1|s) &= c_L \\ p(L = 0|s) &= 1 - c_L \end{aligned}$$

$$\begin{aligned} p(G = 1|s) &= c_G \\ p(G = 0|s) &= 1 - c_G \end{aligned}$$

where  $c_L$  and  $c_G$  are constant values less than 1.

We will use a separate posterior observation function for language and gesture. Combined with the probability to receive a null speech or null gesture observation, the full expression is as follows:

$$\begin{aligned} p(o|s) &= \eta p(l|s)p(g|s) \\ &= \eta p(L|s)p(l|\iota, \beta)^L \cdot p(G|s)p(g|\iota, \beta)^G \end{aligned}$$

*2) Transition Function:* We make the assumption that every component of our state transitions conditionally independent of the other components given the previous state. In addition, each state variable depends only on its previous value and the action taken.

$$\begin{aligned} p(\iota_{t+1}, \beta_{t+1}, \mathcal{I}_{t+1} | \iota_t, \beta_t, \mathcal{I}_t, a_t) \\ = p(\iota_{t+1} | \iota_t, a_t) p(\beta_{t+1} | \beta_t, a) p(\mathcal{I}_{t+1} | \mathcal{I}_t, a_t) \end{aligned}$$

*a) Object Transition Function:* As previously stated, we assume the object the human desires does not change unless its desired object is picked up. If the desired object is picked, it transitions uniformly at random between the remaining objects.

If  $a_t$  is `pick(i)`,

$$p(\iota_{t+1} | \iota_t, a_t) = \begin{cases} 1/(|\mathcal{I}| - 1) & \text{if } \iota_{t+1} \in \mathcal{I} \text{ and } \iota_{t+1} \neq \iota_t \\ 0 & \text{otherwise} \end{cases}$$

If  $a_t$  is any other action:

$$p(\iota_{t+1} | \iota_t, a_t) = \begin{cases} 1 & \text{if } \iota_{t+1} = \iota_t \\ 0 & \text{otherwise} \end{cases}$$

*b) Object Set Transition Function:* The set of objects  $\mathcal{I}$  changes only when an object is picked up. The picked object is removed from the set.

If  $a_t$  is a `pick(i)` action:

$$p(\mathcal{I}_{t+1} | \mathcal{I}_t, a_t) = \begin{cases} 1 & \text{if } \mathcal{I}_{t+1} = \mathcal{I}_t \setminus \{i\} \\ 0 & \text{otherwise} \end{cases}$$

If  $a_t$  is any other action:

$$p(\mathcal{I}_{t+1} | \mathcal{I}_t, a_t) = \begin{cases} 1 & \text{if } \mathcal{I}_t = \mathcal{I}_{t+1} \\ 0 & \text{otherwise} \end{cases}$$

*c) Belief Transition Function:* The human's belief about which object the robot will hand over,  $\beta_t$ , is informed by the actions the robot takes. In the same way that the robot's belief state  $b$  is a summary of all the observations the robot has made [?],  $\beta_t$  is a summary of all observations the human has made of the robot's actions, and reflects the human's estimate about the robot's true state. We therefore update  $\beta_t$  according to Bayesian probabilities in the same way that  $b$  is updated. The following is a standard Markovian update:

For a particular state  $\iota_t$ , its probability after  $t$  observations is given as follows:

$$p(\iota_{t+1} | o_{1:t+1}) = p(o_{t+1} | \iota_{t+1}) \sum_{\iota_t} p(\iota_{t+1} | \iota_t) p(\iota_t | o_{1:t})$$

Recall that  $\beta$  is a distribution over states (objects)  $\iota_t$ , so we can use this rule to write an update for each entry of  $\beta$ .

$$\begin{aligned} \beta_{t+1}(\iota_{t+1}) &= p(\iota_{t+1} | a_{1:t}) \\ &= p(a_t | \iota_{t+1}) \sum_{\iota_t} p(\iota_{t+1} | \iota_t) p(\iota_t | a_{1:t-1}) \\ &= p(a_t | \iota_{t+1}) \sum_{\iota_t} p(\iota_{t+1} | \iota_t) \beta_t(\iota_t) \end{aligned}$$

This requires us to specify  $p(a_t | \iota_{t+1})$ . Recall that actions are operating as observations from the perspective of the human, making this expression a *reverse observation function*. We must provide a reverse observation function for each action. These will often be very similar to the base-level observations described in Section ??, though adapted slightly to suit being performed by a Baxter robot. See section ?? for details.

We must also specify a transition function,  $p(\iota_{t+1} | \iota_t)$ . We will use the object transition function described earlier.

We can rewrite our element-wise update for the whole  $\beta$  vector using a matrix multiplication:

$$\beta_{t+1} = A_t T \beta_t$$

where  $A_t$  is a matrix representing the reverse observation function  $p(a_t | \iota_{t+1})$  and  $T$  is a matrix representing

the transition function  $p(\iota_{t+1}|\iota_t)$ . If the set of items  $\iota = \{\iota_1, \iota_2, \dots, \iota_n\}$ ,  $A$  is a diagonal matrix such that

$$a_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ p(a_t|\iota_{t+1}) = \iota_i & \text{otherwise} \end{cases}$$

and  $T$  is a matrix with entries

$$t_{ij} = p(\iota_{t+1} = \iota_i | \iota_t = \iota_j)$$

This gives us a deterministic transition function for  $\beta$ :

$$p(\beta_{t+1}|\beta_t, a_t) = \begin{cases} 1 & \text{if } \beta_{t+1} = A_t T \beta_t \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

If we assume that  $\beta_0$  begins as the uniform distribution over objects, we can now calculate  $\beta_t$  for any given time step deterministically. However, there are many assumptions that were made, and there is additional expressiveness that can be used by allowing for non-deterministic updates to  $\beta$ . However, we lead the exploration of this space for future work.

*d) Reverse Observation Functions:* We define a reverse observation function for each action the robot can take. These define how  $\beta$  is updated when a robot takes an action.

For `look(x)` actions, we define an angle  $\phi_l$  to be the angle between a vector from the robot's face to the location  $x$  and the vector from the robot's face to the location of object  $\iota_{t+1}$ . We assume this angle is chosen from a normal distribution centered at 0 with a hand tuned variance,  $v_l$ .

$$p(a_t = \text{look}(x) | \iota_{t+1}) = \mathcal{N}(\phi_l | 0, v_l)$$

For `point(x)` actions, we define an angle  $\phi_p$  to be the angle between the vector from Baxter's wrist to location  $x$  and the vector from Baxter's wrist joint to the object  $\iota_{t+1}$ . This gesture is chosen according to a normal distribution with hand-tuned variance  $v_p$ .

$$p(a_t = \text{point}(x) | \iota_{t+1}) = \mathcal{N}(\phi_p | 0, v_p)$$

Actions `say(p)` take the form of a statement "I believe that the object you want is  $p$ ", where  $p$  is some property from some object's vocabulary. The probability of this action given some  $\iota$  is identical to the speech model for the human, given in section ??.

#### D. Policy Generation

**EXW: Will rewrite this depending on whether we use POSS or not. If not, this section will probably removed or changed to say only BSS was used**

It is possible to calculate an optimal policy to solve our POMDP, however, due its size, doing so is intractable. Moreover, we would like to solve our POMDP at speeds suitable for interaction. Fortunately, approximate solutions exist which allow us to trade off accuracy for computational time. One such solver is the Partially Observable Sparse Sampling (POSS) algorithm [? ], which uses Monte-Carlo Tree Search [? ] and Upper Confidence Bound exploration [? ]

] to derive a policy for our POMDP. In order to improve the accuracy of the planner, we also provide a heuristic for the random rollouts that prevents the rollouts from selecting pick actions that are unlikely to select the correct objects.

We implement this model within the Brown UMBC Reinforcement Learning And Planning [? ] framework, which allows for many standard AI and Reinforcement learning algorithms to be used with our domain.

#### IV. EVALUATION

We use Belief Sparse Sampling [? ] (BSS) to solve our domain to evaluate the efficacy of our model. BSS takes several minutes to generate an action, which is too slow for user interaction. Instead, we evaluate the performance of our model in simulation. A procedure for user studies with human participants is also described.<sup>4</sup>

##### A. Results in Simulation

We run 25 trials in simulation. Results are presented in figure ???. We see a reduction in the number of actions needed until the robot can successfully pick the correct object, while accuracy remains the same. With only six objects, it is easy to identify the correct object with enough observations, accounting for the high accuracy in both trials. However, taking social feedback actions allows the agent to request the observations that it needs from the human to complete the interaction as quickly as possible. Notice also that the variance in the number of actions the robot takes until the first correct pick is very small in the social trials, while it is very large in the non-social trials. Social feedback actions allow the agent to control the interaction so that it is as brief as possible.

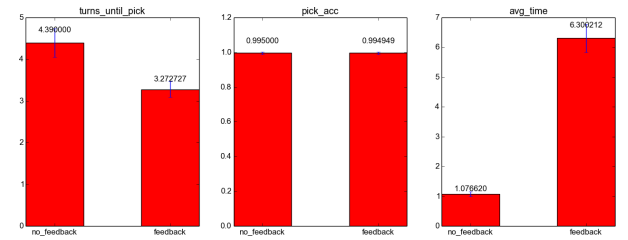


Fig. 5. Results in simulation. `turns_until_pick` is the number of timesteps until first correct pick, `pick_acc` is the ratio of correct picks to total picks, `avg_time` is the time to compute a single action

##### B. User Studies Procedure

**stefie10: cut**

**EXW: Will rewrite this if/when we run user studies**

Users are acquired from convenience sampling by inviting volunteers from the Brown University CIT building. For each user, two trials are run: a baseline trial without any social feedback actions, and a social-feedback trial with social feedback actions enabled. For each trial, the user is

<sup>4</sup>We do note that it would be possible to run the full domain in real time to test with users by precomputing and caching a policy for each belief state that describes the optimal action to take from that state. This is the scheme that we used in [? ]

given a microphone and invited to use speech and gesture to ask the robot for a specific item on the table, and to not change that object until it was picked up. After their requested object is picked up, they request a new object and repeat the procedure until the robot has handed them all objects on the table. It is decided randomly whether the social or non-social trial is performed first for each user. Between the two trials, they are given a survey to assess the qualitative aspects of the interaction, including metrics such as the robot’s perceived intelligence, friendliness, etc.

This procedure is based off a previous user study run for an earlier iteration of this project. In this user study, we noticed a small effect showing the increased speed of interaction with social feedback actions enabled. However, this study was run with only six participants, so these effects were likely not statistically significant. When user studies are run in the future, we will have a larger sample size, but also increase the difficulty of the task by adding more objects, which should further highlight the necessity of social feedback.

Interesting metrics to measure include are how the perception of the robot is correlated to the actions the robot takes, as well as the accuracy and speed of the interaction. We may also consider a longer-running user study to investigate how people learn how to interact with the robot.

## V. FUTURE WORK

The next task will be to run user-studies with human participants to compare the performance of the system with and without social feedback. This will require the use of approximate solvers to allow us to solve the POMDP in real time, or to precompute policies.

We would also like to expand the domain by allowing for more modalities of user input, such as other types of gestures as well as more subtle cues such as expression or emotion.

In the future, we will investigate more sophisticated observation functions, allowing for language interpretation that can understand prepositional phrases and other ways of composing referring expressions.

Other areas for further exploration can be found by removing many of the assumptions made in this thesis. For example, the set of objects on the table was assumed to be known. Attempting to learn from the human more about the set of objects, such as their vocabularies and location, would be an interesting expansion.

Use of the  $\beta$  variable can be expanded to include uncertainty over the value of  $\beta_t$ , allowing for ambiguity in how the human interprets the robot’s action. In the same vein, additional modes of interaction from both the robot and human can be explored.

The object delivery domain can also be abstracted to a domain where the goal is not to determine which object the human wants, but the human’s intent in general. This system of tracking the human’s belief about the robot’s communicated intent can be widely applied to other human-robot interaction tasks.

## VI. ACKNOWLEDGEMENTS

**EXW: I may remove this section, unless there is more to add**

This thesis is the result of work that has been ongoing since December 2014, including a AAAI workshop paper, [? ]. I would like to acknowledge my many collaborators, Advik Guha, Miles Eldon, David Whitney, Horatio Han, Nakul Gopalan, Lawson Wong, James MacGlashan, and John Oberlin.

This project is part of an ongoing initiative towards social robotics by Brown University’s H2R laboratory under Professor Stefanie Tellex.