

An Efficient, Generalized Bellman Update For Cooperative Inverse Reinforcement Learning

Dhruv Malik^{*1} Malayandi Palaniappan^{*1} Jaime F. Fisac¹ Dylan Hadfield-Menell¹ Stuart Russell¹
Anca D. Dragan¹

Abstract

Our goal is for AI systems to correctly identify and act according to their human user’s objectives. Cooperative Inverse Reinforcement Learning (CIRL) formalizes this *value alignment* problem as a two-player game between a human and robot, in which only the human knows the parameters of the reward function: the robot needs to learn them as the interaction unfolds. Previous work showed that CIRL can be solved as a POMDP, but with an action space size exponential in the size of the reward parameter space. In this work, we exploit a specific property of CIRL—the human is a full information agent—to derive an optimality-preserving modification to the standard Bellman update; this reduces the complexity of the problem by an exponential factor and allows us to relax CIRL’s assumption of human rationality. We apply this update to a variety of POMDP solvers and find that it enables us to scale CIRL to non-trivial problems, with larger reward parameter spaces, and larger action spaces for both robot and human. In solutions to these larger problems, the human exhibits pedagogic (teaching) behavior, while the robot interprets it as such and attains higher value for the human.

1. Introduction

As AI agents improve in their ability to optimize for a given objective, it becomes increasingly important that these agents pursue the *right* objective. The *value alignment* problem (Hadfield-Menell et al., 2016; Bostrom, 2014) is that of ensuring that robots optimize for what people want—that robot objectives are aligned with their end-users’ objectives. (We henceforth use robot to refer generically to an AI agent.)

^{*}Equal contribution ¹Department of Electrical Engineering and Computer Sciences, University of California, Berkeley. Correspondence to: Dhruv Malik <dhruvmalik@berkeley.edu>, Malayandi Palaniappan <malayandi@berkeley.edu>.

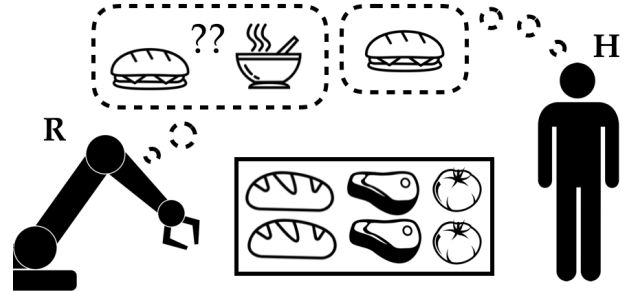


Figure 1. A CIRL game. The human H and the robot R need to work together to prepare a meal. R starts off unaware of which meal H wants, but both H and R get rewarded only if they prepare H ’s desired meal. Solving such a CIRL game has thus far been intractable. In Section 3, we derive a modified Bellman update for computing optimal solutions to CIRL games that achieves an exponential reduction in running time and relaxes CIRL’s assumption of human rationality.

A highly-capable autonomous agent working towards the wrong goal can cause undesired effects, the magnitude of which will tend to increase with the capabilities of the agent. Unfortunately, we humans have a hard time specifying what it is that we actually want. For example, customers may give mistaken instructions to an AI system and system designers may select simple, but potentially incorrect, reward functions to optimize (Amodei & Clark, 2016). Optimizing for the wrong objective can lead to unintended and negative consequences (Amodei et al., 2016).

Rather than optimize a pre-specified reward function, a robot may instead attempt to infer what people *internally* want but cannot perfectly explicate. The robot can use a person’s actions to learn about the reward function over time. The most common approach for this is Inverse Reinforcement Learning (IRL) (Ng & Russell, 2000). IRL makes two implicit assumptions: 1) that the robot is a passive observer, watching the human, and 2) that the human acts as an expert in isolation, ignoring that the robot needs to learn.

Cooperative Inverse Reinforcement Learning (CIRL) (Hadfield-Menell et al., 2016) relaxes these two assumptions. It proposes a formulation in which the human H and the robot R are on the same team and collaborate to achieve the same goal. CIRL is a two-player game between H and

\mathbf{R} , in which *both take actions*, and *both get rewarded according to the same reward function*. The key to CIRL is that only \mathbf{H} knows the parameters of this reward function.

Take for instance the domain from Figure 1. \mathbf{H} and \mathbf{R} work to prepare a meal using three ingredient types: bread, meat, and tomatoes. \mathbf{H} wants to prepare either a sandwich (2 bread, 1 meat, 0 tomatoes), or tomato soup (1 bread, 1 meat, 2 tomatoes). \mathbf{R} does not know *a priori* which meal \mathbf{H} wants, and, to emulate the difficulty that people have in specifying what they want, we assume \mathbf{H} cannot explicate this information directly to \mathbf{R} . At every time step, \mathbf{R} and \mathbf{H} each prepare a single unit of any ingredient, or no ingredient at all. They both receive reward of 1 if they succeed in preparing the right recipe, and 0 otherwise.

In this domain, CIRL captures that the human has an incentive for the robot to infer the correct recipe; and that the robot can take actions in response to the human’s, as opposed to passively waiting until it knows which recipe is right. Crucially, the robot shares the reward function and has an incentive to maximize the human’s internal reward. This creates an incentive to mitigate and avoid unintended consequences from misspecified rewards.

Solving a CIRL game, however, amounts to solving a Dec-POMDP. Previous work has shown that a CIRL game can be reduced to a POMDP. However, the action space in this POMDP is exponential in the size of the reward parameter space. Since POMDP algorithms scale poorly with the size of the action space, non-trivial CIRL games remain difficult to solve with this approach. Additionally, solutions to CIRL are only optimal under the assumption that the human is optimal. This is a strong assumption: it is a well-established fact in cognitive science that humans are often sub-optimal in decision making (Tversky & Kahneman, 1975; Simon, 1957). Our contributions in this paper are three-fold:

1. A Modified Bellman Update: We exploit the fact that the human is a full information agent in CIRL to derive an optimality-preserving modification of the standard Bellman update. This reduces the complexity of the problem by an exponential factor. We show how to apply this modification to existing POMDP solvers (both exact and approximate).

We further show that our modified Bellman update allows us to relax CIRL’s assumption of human rationality. We instead only require that the human’s policy be parameterized by her Q-values. This allows us to solve more realistic CIRL games where the human is modelled as sub-optimal.

2. Empirical Comparison: We show empirically that our method helps scale POMDP solvers to CIRL games with larger reward parameter and action spaces. We find a speed-up of several orders of magnitude for exact methods, and substantial improvements in value for approximate methods.

3. Implications: With the ability to solve more complex CIRL problems, we analyze the solutions that emerge. In

contrast to IRL, we see solutions that exhibit implicit communication. The human takes explicitly suboptimal actions that are better signals for the right reward, and the robot attains higher value for the human because it can take advantage of these signals to learn faster. The coordination that emerges is a consequence of the human and robot being on the same team and reasoning about helping each other.

2. Background

2.1. POMDPs

POMDPs provide a rich model for planning under uncertainty (Sondik, 1971; Kaelbling et al., 1998). Formally, a POMDP is a tuple $\langle X, A, Z, T, O, r, \gamma \rangle$ where X is the set of states; A is the set of the agent’s actions; Z is the set of observations; $T(x_t, a_t, x_{t+1})$ is the transition distribution; $O(x_{t+1}, a_t, z_{t+1})$ is the observation distribution; r is the reward function; and γ is the discount factor.

Consider a simplified instance of the cooking task from Figure 1 where \mathbf{H} picks her actions according to only her desired recipe and the quantity of each ingredient prepared so far, i.e., she does not consider \mathbf{R} ’s past or future behavior when picking her actions. The simplified cooking task is now a POMDP: \mathbf{R} is the agent and \mathbf{H} is a part of the environment. The state specifies \mathbf{H} ’s desired recipe and the quantity of each ingredient already prepared. Thus, \mathbf{H} picks her actions as a function of only the state.

In a POMDP, the agent cannot observe the state; instead, it maintains a belief $b \in \Delta X$, where $b(x)$ is the probability that the agent is in state x . At each time step, the agent receives an observation that helps inform its decisions. The agent in our cooking task, \mathbf{R} , does not know \mathbf{H} ’s desired recipe—a component of the state. \mathbf{R} observes \mathbf{H} ’s actions and tries to infer the desired recipe from \mathbf{H} ’s behavior.

The behavior of the agent is specified by a conditional plan $\sigma = (a, v)$; a denotes the agent’s action and v is a mapping from observations to future conditional plans for the agent to follow. An example conditional plan for \mathbf{R} is: prepare meat now and if \mathbf{H} responds by preparing bread, prepare a second slice of bread; if \mathbf{H} prepares tomatoes, prepare another batch of tomatoes; or if \mathbf{H} prepares meat, do not prepare any ingredient.

The α -vector of a conditional plan contains the value of following the plan at any given state:

$$\alpha_\sigma(x) = R(x) + \gamma \sum_{x' \in X} \sum_{z \in Z} P(x', z | x, a) \alpha_{v(z)}(x') \quad (1)$$

The value of a plan at a belief b is the expected value of the plan across the states i.e. $V_\sigma(b) = b \cdot \alpha_\sigma = \sum_{x \in X} b(x) \alpha_\sigma(x)$. The goal of an agent in a POMDP is to find the plan with maximal value from her current belief.

Value iteration (Sondik, 1971) can be used to compute the optimal conditional plan. This algorithm starts at the horizon

and works backwards. It generates new conditional plans at each time-step and evaluates them according to Eq. 1. It constructs all potentially optimal plans of length T and selects the one with maximal value at the initial belief.

2.2. Cooperative Inverse Reinforcement Learning

Now, consider an instance of the cooking task where \mathbf{H} is a second agent in the game and no longer behaves independently of \mathbf{R} . There is now a strong interdependence between \mathbf{H} 's and \mathbf{R} 's behavior: \mathbf{H} 's actions both depend on and influence \mathbf{R} 's belief. This problem is now no longer a POMDP; it is a CIRL game.

A CIRL game is an asymmetric-information two-player game between a human \mathbf{H} and a robot \mathbf{R} (Hadfield-Menell et al., 2016). \mathbf{H} knows the true reward function and \mathbf{R} does not initially. Formally, a CIRL game is a tuple: $M = \langle X, \{\mathcal{A}^H, \mathcal{A}^R\}, T, \{\Theta, r\}, \gamma \rangle$. X is the set of observable world-states; \mathcal{A}^H and \mathcal{A}^R are the actions available to \mathbf{H} and \mathbf{R} respectively; $T(x_t, a_t^H, a_t^R, x_{t+1})$ is the transition distribution; Θ is the set of reward parameters; r is the parameterized reward function shared by both agents; γ is the discount factor. A solution to a CIRL game is a pair of policies—one for \mathbf{H} and \mathbf{R} each—that maximizes the expected reward obtained by \mathbf{H} and \mathbf{R} .

In our cooking task, Θ is the set of possible recipes. \mathbf{R} does not know \mathbf{H} 's desired recipe, $\theta \in \Theta$. The reward function r is parameterized by Θ : both agents receive a reward of 1 if they succeed in preparing \mathbf{H} 's desired recipe.

Reducing a CIRL game to a POMDP A CIRL game is a Dec-POMDP (Bernstein et al., 2002) but it can be reduced to a POMDP where the optimal policy corresponds to optimal CIRL policy pairs (Hadfield-Menell et al., 2016). The states in this POMDP are tuples of world-state and reward parameter: $S = X \times \Theta$; the actions are tuples (δ^H, a^R) specifying a decision rule $\delta^H : \Theta \rightarrow \mathcal{A}^H$ for \mathbf{H} , which maps reward parameters to human actions, and an action a^R for \mathbf{R} ; the observations are \mathbf{H} 's action at the last time step.

An example action in the reduced POMDP of the cooking task is a tuple, where the first entry specifies that \mathbf{H} prepares bread if she prefers a sandwich and prepares tomatoes if she prefers soup, and the second entry specifies that \mathbf{R} prepares bread (regardless of its belief).

This reduction enables us to solve a CIRL game using POMDP algorithms. However, the size of the action space in this POMDP is $|\mathcal{A}^H|^{|\Theta|} |\mathcal{A}^R|$, as shown in Figure 2. (There are $|\mathcal{A}^H|^{|\Theta|}$ possible decision rules for \mathbf{H} and $|\mathcal{A}^R|$ actions for \mathbf{R} .) In other words, the action space in this POMDP grows exponentially with the size of the reward parameter space. Exact POMDP algorithms are exponential in the size of the action space, so this approach can only be applied to very small CIRL problems.

Additionally, the policy for \mathbf{R} that is output by the reduced

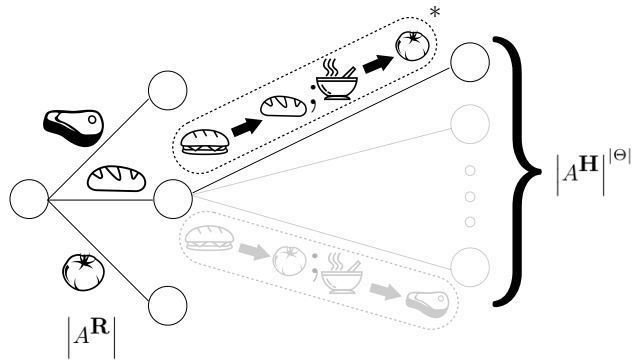


Figure 2. A node in the search tree from the POMDP reduction of our example CIRL game. Actions are tuples that contain an action for \mathbf{R} and a decision rule for \mathbf{H} – a mapping from her desired recipe to an action. This leads to a branching factor of $|\mathcal{A}^H|^{|\Theta|} |\mathcal{A}^R|$ and makes application of POMDP methods inefficient. In Section 3.2, we derive a modified Bellman update that prunes away all of \mathbf{H} 's decision rules but the optimal response. (In the diagram, the gray branches are pruned away by our modified Bellman update.)

POMDP is optimal only if \mathbf{H} is perfectly rational, i.e., if \mathbf{H} is guaranteed to always pick the optimal action. This is an unrealistic assumption: humans are not idealized rational agents (Tversky & Kahneman, 1975; Simon, 1957).

3. A Modified Bellman Update for CIRL

If \mathbf{H} were following a fixed policy based on the state $s = (x, \theta)$, we could encode \mathbf{H} as a part of the environment. However, in the interactive setting of a CIRL game, \mathbf{H} may plan for changes in \mathbf{R} 's belief. If we encode \mathbf{H} in the environment, the dynamics change in response to \mathbf{R} 's intended plan and the problem is no longer a POMDP. Our main contribution is to derive a modified Bellman update for POMDP algorithms to solve this problem.

Our key idea is as follows. During planning, we know \mathbf{R} 's intended future response to each of \mathbf{H} 's actions. \mathbf{H} has full state information, so the α -vectors in value iteration allow us to directly compute \mathbf{H} 's Q-values. We can therefore also compute her optimal action based on \mathbf{R} 's intended future response. This means we do not have to reason over the set of decision rules for \mathbf{H} : we can solve a CIRL game by instead solving a POMDP with time-varying dynamics and, importantly, where the action space has size $|\mathcal{A}^R|$. This is exponentially smaller than the action space of size $|\mathcal{A}^H|^{|\Theta|} |\mathcal{A}^R|$ in the reduced POMDP. This amounts to a modified Bellman update.

3.1. The Transition Dynamics

If \mathbf{H} follows a policy that depends only on the state $s = (x, \theta)$, the dynamics of the game can be computed as:

$$\begin{aligned}
 P(s', a^H | s, a^R) &= P((x', \theta'), a^H | (x, \theta), a^R) \\
 &= P((x', \theta') | (x, \theta), a^H, a^R) \cdot P(a^H | (x, \theta), a^R) \\
 &= T(x, a^H, a^R, x') \cdot \mathbf{1}(\theta = \theta') \cdot P(a^H | x, a^R, \theta) \\
 &\stackrel{a.}{=} T(s, a^H, a^R, s') \cdot P(a^H | x, a^R, \theta)
 \end{aligned}$$

However, in the CIRL formulation, \mathbf{H} does not behave according to a fixed policy. \mathbf{H} , who is assumed to be rational, behaves according to her Q-values and picks the action that maximizes her expected value. Due to the interdependence between \mathbf{H} 's and \mathbf{R} 's behavior, these Q-values depend on \mathbf{R} 's conditional plan. The dynamics then are:

$$\begin{aligned}
 P(s', a^H | s, \sigma) &= P((x', \theta'), a^H | (x, \theta), (a^R, v)) \\
 &= T(x, a^H, a^R, x') \cdot \mathbf{1}(\theta' = \theta) \cdot P(a^H | x, a^R, v, \theta) \\
 &= T(s, a^H, a^R, s') \cdot \mathbf{1}(a^H = \arg \max_{a^H} Q_H(s, a^H, \sigma))
 \end{aligned} \tag{2}$$

These dynamics change over time since they depend on the robot's future behavior. However, \mathbf{R} 's behavior depends on these dynamics, so, we cannot pre-compute them as part of a POMDP reduction. However, we do have access to \mathbf{R} 's future conditional plan *during* planning. This means we can compute \mathbf{H} 's Q-values, and, consequently, the transition probabilities, with a modification to the Bellman update.

3.2. Adapting POMDP Value Iteration

POMDP value iteration rolls back the values of the game from the horizon, storing them as α -vectors. If \mathbf{R} follows a plan $\sigma = (a^R, v)$, then we can compute \mathbf{H} 's Q-values as

$$Q_H(s, a^H, \sigma) = \sum_{s'} T(s, a^H, a^R, s') \cdot \alpha_{v(a^H)}(s').$$

\mathbf{H} 's optimal action maximizes this expression. To leverage this, we adapt the Bellman update in Eq. 1 to replace the dynamics of the game with $P(s', a^H | s, \sigma)$ from Eq. 2. The modified Bellman update is then:

$$\alpha_\sigma(s) = R(s) + \gamma \cdot \max_{a^H} \sum_{s' \in S} T(s, a^H, a^R, s') \cdot \alpha_{v(a^H)}(s'). \tag{3}$$

We can then use value iteration with this modified Bellman update to compute the \mathbf{R} 's optimal policy. The following theorem establishes that the modified Bellman update, Eq. 3, is optimality-preserving.

Theorem 1. *For any CIRL game, the policy computed by value iteration with the modified Bellman update is optimal.*

All theorem proofs are presented in Appendix A.

a. We let $T(s, a^H, a^R, s') = T(x, a^H, a^R, x') \cdot \mathbf{1}(\theta = \theta')$.

This modification to the Bellman update allows us to solve a CIRL game without having to include the set of \mathbf{H} 's decision rules in the action space. As depicted in Figure 2, the modified Bellman update computes \mathbf{H} 's optimal action given the current state and the robot's plan; all of \mathbf{H} 's other actions are pruned away in the search tree. The size of the action space is then $|\mathcal{A}^R|$ instead of $|\mathcal{A}^H|^{|\Theta|} |\mathcal{A}^R|$. POMDP algorithms are exponential in the size of the action space; this modification therefore allows us to solve CIRL games much more efficiently. The following theorem establishes the complexity gains made by algorithm.

Theorem 2. *The modification to the Bellman update presented above reduces the time and space complexity of a single step of value iteration by a factor of $\mathcal{O}(|\mathcal{A}^H|^{|\Theta|})$.*

3.3. Relaxing CIRL's Assumption of Rationality

To achieve value alignment, we can now efficiently solve a CIRL game to find an optimal policy for \mathbf{R} . However, this policy is optimal only if \mathbf{H} is perfectly rational: a strong assumption. This is rarely true in reality; we thus want to find an optimal policy for \mathbf{R} even when \mathbf{H} is sub-optimal.

In addition to improving efficiency, our modified Bellman update allows us to do exactly that and relax CIRL's assumption of rationality. The dynamics of our modified Bellman update, presented above as Eq. 2, do not require that \mathbf{H} is perfectly rational. These dynamics will remain well-defined so long as we know the distribution over \mathbf{H} 's actions, π_H , and can compute the probability that she picks any action from her current state. To avoid compromising the interactive nature of CIRL, we require that π_H must be a function of \mathbf{H} 's Q-values, which account for the robot's future behavior. The dynamics of the game are then given by:

$$P(s', a^H | s, \sigma) = T(s, a^H, a^R, s') \cdot \pi_H(a^H | Q_H(s, a^H, \sigma)).$$

The modified Bellman update is then:

$$\begin{aligned}
 \alpha_\sigma(s) &= R(s) + \gamma \cdot \sum_{a^H} \pi_H(a^H | Q_H(s, a^H, \sigma)) \cdot \\
 &\quad \sum_{s' \in S} T(s, a^H, a^R, s') \cdot \alpha_{v(a^H)}(s').
 \end{aligned} \tag{4}$$

With this modified Bellman update, we may now use value iteration to solve CIRL games without assuming that the human is perfectly rational. We instead only require that the human selects her actions with respect to her Q-values. This restriction is rather broad and includes a variety of models of human decision making from cognitive science. A popular example of such a model is Boltzmann-rationality, where the human picks her actions according to a Boltzmann distribution over her Q-values, i.e.,

$$\pi_H(a^H | Q_H(s, a^H, \sigma)) \propto \exp(\beta \cdot Q_H(s, a^H, \sigma))$$

where β is a parameter which controls how rational the human is. (A higher β corresponds to a more rational human.)

Algorithm 1 Adapted Value Iteration for CIRL Games

```

1:  $\Gamma_t \leftarrow$  Set of trivial plans
2: for  $t \in \{T-1, T-2, \dots, 1, 0\}$  do
3:    $\Gamma_{t+1} \leftarrow \Gamma_t$ 
4:    $\Gamma_t \leftarrow$  Set of all plans beginning at time  $t$ 
5:   for  $\sigma \in \Gamma_t$  do
6:     for  $s = (x, \theta) \in S$  do
7:        $Q_H(s, a^H, \sigma) = \sum_{s'} T(s, a^H, a^R, s')$ 
8:        $\alpha_{v(a^H)}(s')$ 
9:        $\alpha_\sigma(s) = R(s) + \gamma \cdot \sum_{a^H} \pi_H(a^H | Q_H(s, a^H, \sigma)) \cdot Q_H(s, a^H, \sigma)$ 
10:       $\pi_H(a^H | Q_H(s, a^H, \sigma)) \cdot Q_H(s, a^H, \sigma)$ 
11:     end for
12:   end for
13:    $\Gamma_t \leftarrow$  Prune( $\Gamma_t$ )
14: end for
15:  $a_*^R = \operatorname{argmax}_{\sigma \in \Gamma_0} \alpha_\sigma \cdot b_0$ 
16: Return  $a_*^R$ 

```

The time and space complexity of value iteration with this Bellman update is identical to that with the modified Bellman update presented in Section 3.2, and analyzed in Theorem 2. The pseudocode for our adapted algorithm is presented as Algorithm 1 below.

4. Adapting Approximate Algorithms

Approximate algorithms for POMDPs often rely on variants of the Bellman update. This lets us use our modified Bellman update to improve approximate algorithms for CIRL.

4.1. PBVI

Background Point Based Value Iteration (PBVI) is an approximate algorithm used to solve POMDPs (Pineau et al., 2003). The algorithm maintains a representative set of points in belief space and an α -vector at each of these belief points. It performs approximate value backups at each of these belief points using this set of α -vectors. Let Γ_{t+1} denote the set of α -vectors for plans that begin at time $t+1$. The value at time t at a belief b is approximated as:

$$V(b) = \max_{a \in A} \left[\sum_{s \in S} R(s) b(s) + \gamma \sum_{o \in O} \max_{\alpha \in \Gamma_{t+1}} \sum_{s \in S} \left(\sum_{s' \in S} P(s', o | s, a) \alpha(s) \right) b(s) \right].$$

The algorithm trades off computation time and solution quality by expanding the set of belief points over time: it randomly simulates forward trajectories in the POMDP to produce new, reachable beliefs.

Our Adaptation If \mathbf{R} takes action a^R and follows a conditional plan σ , then \mathbf{H} 's Q-values are $Q_{\mathbf{H}}(x, a^H, a^R, \alpha) = \sum_{s'} T(s, a^H, a^R, s') \cdot \alpha_\sigma(s')$. Notice that we can compute these Q-values at each step of PBVI. This lets us use the

modified Bellman update and to adapt PBVI to solve CIRL games specifically. We replace the transition-observation distribution in the PBVI backup rule with

$$P(s', a^H | s, a^R, \alpha) = T(s, a^H, a^R, s') \cdot \pi_{\mathbf{H}}(Q_{\mathbf{H}}(x, a^H, a^R, \alpha)).$$

The modified backup rule for PBVI is thus given by

$$V(b) = \max_{a^R \in \mathcal{A}^R} \left[\sum_{s \in S} R(s) b(s) + \gamma \sum_{a^H} \max_{\alpha \in \Gamma_{t+1}} \sum_{s \in S} \left(\sum_{s' \in S} P(s', a^H | s, a^R, \alpha) \alpha(s) \right) b(s) \right].$$

We now show that the approximate value function in PBVI converges to the true value function. Let ϵ_B denote the density of the set of belief points B in PBVI. Formally, $\epsilon_B = \max_{b \in \Delta} \min_{b' \in B} \|b - b'\|_1$ is the maximum distance from any reachable, legal belief to the set B .

Theorem 3. For any belief set B and horizon n , the error of our adapted PBVI algorithm $\eta = \|V_n - V_n^*\|_\infty$ is bounded as

$$\eta \leq \frac{(R_{\max} - R_{\min}) \epsilon_B}{(1 - \gamma)^2}.$$

4.2. POMCP

Background POMCP is a Monte Carlo tree-search (MCTS) based approximate algorithm for solving large POMDPs (Silver & Veness, 2010). The algorithm constructs a search tree of action-observation histories and uses Monte Carlo simulations to estimate the value of each node in the tree. During search, actions within the tree are selected by UCB1. This maintains a balance between exploiting actions known to have good return and exploring actions not yet taken (Kocsis & Szepesvári, 2006). At leaf nodes, a rollout policy accrues reward which is then backed up through the tree. The algorithm estimates the belief at each node by keeping track of the hidden state from previous rollouts.

POMCP scales well with the size of the state space, but not with the size of the action space, which determines the branching factor in the search tree. POMCP is thus ill-suited to solving the reduced POMDP of CIRL games since the size of the action space is $|\mathcal{A}^H|^{|\Theta|} |\mathcal{A}^R|$.

Our Adaptation Using the idea behind our modified Bellman update, we adapt POMCP to solve CIRL games more efficiently. We approximate \mathbf{H} 's policy while running the algorithm (much like we exactly compute \mathbf{H} 's policy in exact value iteration). We maintain a live estimate of the sampled Q-values for \mathbf{H} at each node. With enough exploration of the search tree (for instance, if actions are selected using UCB1), the estimated Q-values converge to the true values (in the limit). This guarantees that \mathbf{H} 's policy converges to her true policy. The following result establishes convergence of our algorithm.

Theorem 4. *With suitable exploration, the value function constructed by our adapted POMCP algorithm converges in probability to the optimal value function, $V(h) \rightarrow V^*(h)$. As the number of visits $N(h)$ approaches infinity, the bias of the value function $\mathbb{E}[V(h) - V^*(h)]$ is $O(\log(N(h))/N(h))$.*

The pseudocode for our adapted PBVI and our adapted POMCP algorithm is presented as Algorithm 1 and 2 respectively in Appendix B.

5. Related Work

POMDP Algorithms We chose to explicate our modified Bellman update in the context of PBVI and POMCP because they are the seminal point-based and MCTS algorithms respectively, for solving POMDPs. For example, SARSOP (Kurniawati et al., 2008) and DESPOT (Ye et al., 2017), two state-of-the-art algorithm for POMDPs, are variants of PBVI and POMCP respectively. The principles we outlined in Sections 3 and 4 can be easily adapted to a large variety of point-based and MCTS algorithms, including any which may be developed in the future, to derive even more efficient algorithms for solving CIRL games.

MOMDP Algorithms The POMDP representation of CIRL is also a mixed-observability Markov decision process (MOMDP) since the state space can be factored into a fully- and a partially-observable component. This structure allows for more efficient solution methods; Ong et al. (2010) leverage the factored nature of the state space to create a lower dimensional representation of belief space. This core idea is orthogonal to ours, which exploits CIRL’s information asymmetry instead. The two can be leveraged together.

Dec-POMDP Algorithms Dec-POMDP algorithms can be used to solve CIRL directly, without using the POMDP reduction. These solution methods are generally intractable, but recent work has made progress on this front. Such Dec-POMDP algorithms which attempt to prune away unreasonable strategies resemble our approach. Amato et al. (2009) use reachability analysis to identify reachable states, then consider all policies which are useful at such states. Hansen (2004) model other agents possible strategies as part of a players belief, and prune away weakly dominated strategies at each step. While such approaches use heuristics to prune away some suboptimal strategies, we leverage the information structure of CIRL to compute the optimal strategy for \mathbf{H} and prune away *all* other strategies. This guarantees an exponential reduction in complexity while preserving optimality; this is not true for the other methods.

Value Alignment Recent work has explored relaxing the rationality requirement of CIRL (Fisac et al., 2017). Our work improves on their relaxation in several ways: (1) Their Bellman update assumes that the human acts Boltzmann-rationally. Our modification can model a large variety of

Table 1. Time taken (s) to find the optimal robot policy using exact VI and our adaptation of it for various numbers of possible recipes. NA denotes that the algorithm failed to solve the problem.

# RECIPES	EXACT VI	OURS
2	4.448 ± 0.057	0.071 ± 0.013
3	394.546 ± 6.396	0.111 ± 0.013
4	NA	0.158 ± 0.003
5	NA	0.219 ± 0.007
6	NA	0.307 ± 0.005

human behaviors (including this). (2) Their discretized belief value iteration algorithm has neither the guarantee of optimality of our adapted VI algorithm nor the scalability of our adapted POMCP/PBVI algorithms.

6. Experiments

We now verify that our modified Bellman update allows POMDP algorithms to solve CIRL games more efficiently than the standard update. We ran three experiments: one for exact value iteration (VI), PBVI, and POMCP each. The results of the PBVI experiment are presented in Appendix C.1 due to space constraints. To verify the results of our experiments, we ran further experiments on a second domain. The details and results of these experiments are presented in Appendix C.2.

6.1. Experimental Design

Domain Our experimental domain is based on our running example from Section 1. Assume there are m recipes and n ingredients. The state space is an n -tuple representing the quantity of each ingredient prepared thus far. At each time step, each agent can prepare any of the n ingredients or none at all. Each of the m recipes corresponds to a different θ (i.e. reward parameter) value. Both agents receive a reward of 1 if \mathbf{H} ’s desired recipe is prepared correctly and a reward of 0 otherwise. The robot \mathbf{R} begins the game entirely uncertain about \mathbf{H} ’s desired recipe i.e. \mathbf{R} has a uniform belief over Θ .

We want to stress that this experimental domain is not trivial: one of the domains we managed to solve in our experiments had $\sim 10^{10}$ states.

Manipulated Variables Our primary variable is the type of Bellman update used: modified vs. standard. We also varied the number of recipes, i.e., size of the reward parameter space, and the number of ingredients, i.e., size of \mathbf{H} ’s and \mathbf{R} ’s action space.

Dependent Measures In our first experiment (exact VI), we measured the time taken by the algorithms to solve the problem. In our second experiment (PBVI), we measured the value attained by the algorithms in a fixed time. In our third experiment (POMCP), we measured the value attained by the algorithms in a fixed number of samples.

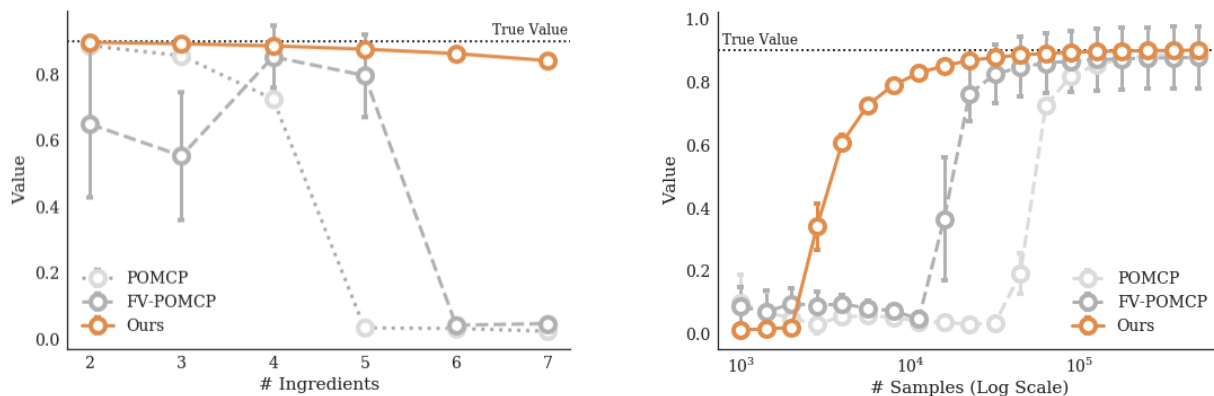


Figure 3. (Left) The value attained by POMCP, FV-POMCP, and our adapted algorithm in 30,000 samples with various numbers of ingredients. (Right) Value attained by POMCP, FV-POMCP and our approximate algorithm with 2 recipes and 6 ingredients.

Hypothesis POMDP algorithms are more efficient at solving CIRL games with the modified Bellman update than with the standard Bellman update.

6.2. Analysis

Exact VI In our first experiment, we compared the time taken by exact VI and by our adaptation of it with the modified Bellman update. We first fixed the number of ingredients at two and varied the number of recipes in the domain. Table 1 compares the results. For the simpler problems, where the number of recipes was 2 or 3, our adapted algorithm solved the problem up to $\sim 3500\times$ faster than exact VI. On more complex problems where the number of recipes is greater than 3, exact VI failed to solve the problem after depleting our system’s 16GB memory; in contrast, our adapted algorithm solved each of these more complex problems in less than 0.5 seconds. We next fixed the number of recipes and compared the performance of both these algorithms for various numbers of ingredients. Both the exact methods, but especially the one using the standard Bellman update, scaled much worse with the number of ingredients than with the number of recipes. With even three ingredients, exact VI timed out and failed to solve the problem within two hours; our algorithm however solved the problem in five seconds.

POMCP We compared the value attained in 30,000 samples by POMCP and by our adaptation with the modified Bellman update. We additionally compared these algorithms with FV-POMCP, a state-of-the-art MCTS method for solving MPOMDPs, a type of Dec-POMDP in which all agents can observe each others’ behavior (as in CIRL).

We first fixed the number of recipes at 2 and varied the number of ingredients. Our adapted algorithm outperformed the other two algorithms across the board, especially for large numbers of ingredients. The results of this comparison are presented in Figure 3 (left). POMCP did poorly on games with more than 4 ingredients. Although FV-POMCP scaled

better to more complex games than POMCP, its values had high variance. For the largest games, with 6 and 7 ingredients, our adapted algorithm was the only one capable of solving the problem in 30,000 iterations. We also compared the value attained by each algorithm across 500,000 samples on the 6 ingredient game. The results of this comparison are depicted in Figure 3 (right). Our algorithm converged to the true value faster than either of the other algorithms.

We next fixed the number of ingredients at 4 and varied the number of recipes. We found that the results of this experiment broadly matched the results of our previous experiment where we varied the number of ingredients. For example, with 4 recipes, our method achieves an average value of 0.631 ± 0.221 in 30,000 iterations while POMCP gets 0.429 ± 0.183 and FV-POMCP gets 0.511 ± 0.124 .

These results together demonstrate that POMDP algorithms with our modified Bellman update scales much better to more complex CIRL games than with the standard Bellman update. This offers strong evidence for our hypothesis.

7. Discussion

The previous section showed that we can now solve larger, non-trivial CIRL games. While we are still far from addressing value alignment in the high dimensional and continuous real world, our work allows us to analyze CIRL solutions to non-trivial problems and understand their implications for value alignment.

7.1. CIRL vs IRL

In the absence of CIRL solutions, a standard approach to learning the human’s internal reward is Inverse Reinforcement Learning (IRL) (Ng & Russell, 2000). We thus compare what advantages CIRL has compared to IRL. On a collaborative task, IRL is equivalent to assuming that \mathbf{H} chooses her actions in isolation, and \mathbf{R} uses observations of

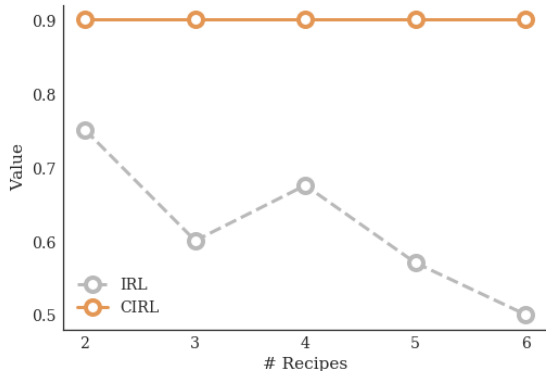


Figure 4. Value attained by CIRL and standard IRL on the cooking domain with various numbers of possible recipes. Unlike IRL, CIRL produces solutions where \mathbf{H} picks her actions pedagogically and \mathbf{R} reasons about \mathbf{H} accordingly.

her behavior to infer her preferences. Specifically, \mathbf{H} solves a single-agent, fully-observable, variant of the problem, and \mathbf{R} responds by solving the POMDP described in Section 2.

We fix the number of ingredients at 3 and vary the number of recipes. Figure 4 shows the results. In each experiment the optimal CIRL solution prepares the correct recipe while the IRL solution fails to do so up to 50% of the time. To understand the nature of this difference in performance, we analyze the CIRL and IRL solutions. Consider a case of our running example with two recipes. The state is a tuple $(\#meat, \#bread, \#tomatoes)$ and $\Theta = \{sandwich = (1, 2, 0), soup = (1, 1, 2)\}$. For both approaches, \mathbf{R} initially prepares meat. In the baseline IRL solution, \mathbf{H} can initially make any ingredient if she wants soup and can make meat or bread if she wants a sandwich. In each case, she chooses uniformly at random between allowed ingredients. This conveys some information about her desired recipe, but is not enough to uniquely identify it. Since the same ingredient is optimal for multiple recipes, \mathbf{R} is still confused after one turn. This means \mathbf{R} will sometimes fail to complete the desired recipe, reducing average utility.

The CIRL solution, in contrast, relies on the implicit communication between the human and the robot. Here, if \mathbf{H} wants soup, she prepares tomatoes, as opposed to any ingredients that are common with the sandwich. Even more interestingly, she waits (i.e. does nothing) if she wants a sandwich. This is pure signaling behavior—*waiting is suboptimal in isolation*, but picking an ingredient is more likely to confuse the robot. In turn \mathbf{R} knows that \mathbf{H} would have picked tomatoes if she wanted soup, and responds appropriately.

In other words, \mathbf{H} teaches the robot about her preferences with her action selection. This works because \mathbf{H} knows that \mathbf{R} will interpret her behavior pragmatically, i.e., \mathbf{R} expects to be taught by \mathbf{H} . This is reflected in the experiment: the optimal CIRL solution prepares the correct recipe each time.

The value alignment problem is necessarily cooperative: without the robot, the human is unable to complete her desired task, and without explicit signaling from the human, the robot learns inefficiently, is less valuable and is more likely to make a mistake. Pedagogic behavior from \mathbf{H} naturally falls out of the CIRL solution. In response, \mathbf{R} interprets \mathbf{H} 's actions pragmatically. These instructive and communicative behaviors allow for faster learning and create an opportunity to generate higher value for the human.

7.2. CIRL with Suboptimal Humans

To further investigate the performance of CIRL in realistic settings (e.g., where \mathbf{H} may not be rational), we ran another experiment. We varied whether \mathbf{H} behaved according to CIRL or IRL, \mathbf{R} 's model of \mathbf{H} in training (rational or Boltzmann-rational), and the actual model of \mathbf{H} (same as previous). We measured the proportion of times they prepared the correct recipe in each setting, fixing the number of ingredients at 3 and recipes at 4. Figure 6 in Appendix D shows the results. (We also conducted a more comprehensive experiment with 20 human behaviors instead of 2. Results are presented in Appendix E.)

Averaged across different models of \mathbf{H} used to train \mathbf{R} , when \mathbf{H} behaved according to CIRL, \mathbf{H} and \mathbf{R} succeeded in preparing the correct recipe $> 90\%$ of the time. This was also true when \mathbf{H} behaved Boltzmann-rationally. This suggests that the pedagogic behavior that arises from CIRL makes it more robust to any sub-optimality from \mathbf{H} . In contrast, when \mathbf{H} behaved as in IRL (i.e., not pedagogically), they only prepared the correct recipe $\sim 70\%$ of the time when \mathbf{H} was rational, and $\sim 40\%$ of the time when \mathbf{H} was not. So, the importance of pedagogic behavior from \mathbf{H} to achieve value alignment is clear.

7.3. Do People Behave Pedagogically?

A question arises at this point as to whether *real* people will adopt the pedagogic behavior predicted by CIRL solutions. To start testing this, we ran a (very preliminary) pilot study to start investigating whether CIRL improves interactions with real people. The details of this experiment are presented in Appendix F. We found some encouraging evidence that suggests people do indeed behave pedagogically when collaborating with a robot; and that a CIRL-trained robot is can be better at exploiting this pedagogic behavior to achieve fluid human-robot collaboration than an IRL-trained robot.

In future work, we plan to conduct a more extensive human subjects study to validate these preliminary findings. We additionally plan to explore techniques to make the robot better elicit pedagogic behavior from the human in their interaction and to make CIRL robust to variations in human behavior.

Appendix

A. Proofs

In this section, we present the proofs for the propositions and theorems in the main paper.

Theorem 1. For any CIRL game, the policy computed by value iteration with the modified Bellman update is optimal.

Proof. During POMDP value iteration, values are propagated from the horizon in the form of α vectors, which store the expected values of executing a conditional plan from each state. In CIRL games, \mathbf{H} is a full information actor, which means she has knowledge of her true reward parameter θ . So, given \mathbf{R} 's conditional plan, she can directly compute her Q values from the α vectors. The presence of the inner max in our modified update rule ensures that \mathbf{H} always picks the action corresponding to her maximum Q value. This implies that at every backup step, our modified Bellman update never prunes away the optimal action. Hence, the output policy, which is the best policy among those not pruned away, must be optimal. \square

Theorem 2. The modification to the Bellman update presented above reduces the time and space complexity of a single step of value iteration by a factor of $\mathcal{O}(|\mathcal{A}^H|^{\Theta})$.

Proof. The complexity of one step of POMDP value iteration is linear in the size of the action space, $|A|$ (Russell & Norvig, 1995). Since the structure of our algorithm is identical to that of exact value iteration, this is also true for our algorithm.

The action space in the POMDP reduction of CIRL has size $|\mathcal{A}^H|^{\Theta}|\mathcal{A}^R|$. Our modification to the Bellman update reduces the size of the action space to simply $|\mathcal{A}^R|$. Therefore, our algorithm reduces the time taken to run, and number of plans generated at, each time step by a factor of $|\mathcal{A}^H|^{\Theta}$. \square

Theorem 3. For any belief set B and horizon n , the error of our adapted PBVI algorithm $\eta = \|V_n - V_n^*\|_\infty$ is bounded as

$$\eta \leq \frac{(R_{\max} - R_{\min})\epsilon_B}{(1 - \gamma)^2}.$$

Proof. Since the dynamics of our problem are now time-varying instead of static, the backup operator applied at every time-step changes. Let H_t denote the backup operator applied to compute the value of V_t . It will suffice to show that each such backup operator H_t is a contraction mapping. The result then follows by following the proof of Theorem 1 in (Pineau et al., 2003) exactly. We will prove the result by

showing that each H_t is the backup operator for a specific POMDP and thus, for this POMDP's corresponding belief MDP; it must therefore be a contraction mapping.

Take H_t for some $1 \leq t \leq n$. We will now construct a new POMDP for which H_t is the backup operator. Let $\hat{S} = S \times \Gamma_{t+1}$, where Γ_{t+1} denotes the set of α -vectors from our original problem at time $t + 1$. Let the α -vector component of the state be static i.e. $P((s', \alpha') \mid (s, \alpha), a^R) = 0$ if $\alpha \neq \alpha'$. The action and observation spaces remain as they are in our original problem. The transition-observation distribution, given in Eqn. (3), is now time-invariant: we do not need to look forward in the search tree to compute the Q-values since the α -vectors are available in the state space. Hence, this POMDP is well-defined.

The dynamics of the POMDP are static and identical to the dynamics of our problem at time t . Thus, the backup operator H_t is the backup operator for this POMDP and also for this POMDP's corresponding belief MDP. Therefore, the backup operator H_t must be a contraction mapping. \square

Theorem 4. With suitable exploration, the value function constructed by our adapted POMCP algorithm converges in probability to the optimal value function, $V(h) \rightarrow V^*(h)$. As the number of visits $N(h)$ approaches infinity, the bias of the value function $\mathbb{E}[V(h) - V^*(h)]$ is $\mathcal{O}(\log(N(h))/N(h))$.

Proof. We will show that with enough exploration, in the limit of infinite samples, we have a well defined POMDP. The result then follows from Theorem 1 in (Silver & Veness, 2010).

The human action nodes in the search tree maintain an array of values, which store the values of picking that action for different θ . At any point in the search tree, the human actions (like the robot actions) are selected by picking the one that has the maximum augmented value (current estimated value plus exploration bonus). So, in the limit of infinite samples, each human action node is visited infinitely many times and the value estimates at the nodes converge to the true Q values. Having the correct human Q values gives us a POMDP, with well defined transition-observation dynamics. The result then follows from applying the analysis given in Theorem 1 of (Silver & Veness, 2010) to this POMDP. \square

B. Pseudocode

The pseudocode for adapted PBVI is presented below. The algorithm follows a similar structure to the standard PBVI algorithm (Pineau et al., 2003). The main difference between our adapted algorithm and the standard PBVI algorithm is that ours uses the modified Bellman update instead of the standard one. See lines 15-16.

The pseudocode for adapted POMCP is also presented below, similar in style to that presented in (Silver & Veness, 2010). The key difference is that we maintain a live estimate of the sampled Q-values for \mathbf{H} at each node. We maintain arrays which store the estimated values of taking each human action for different θ . The optimal human action is selected with regard to these estimates. Like the robot actions, the human actions are selected while balancing exploration and exploitation. Rollouts use random action selection.

Algorithm 2 Adapted PBVI for CIRL Games

```

1: function PBVI ( $b_0, T$ )
2:    $B \leftarrow \{b_0\}$ 
3:    $V \leftarrow$  Set of  $\alpha$ -vectors belonging to trivial plans
4:   repeat
5:     for  $t \in \{T-1, T-2, \dots, 1, 0\}$  do
6:        $V \leftarrow$  Backup( $B, V$ )
7:     end for
8:      $B \leftarrow$  Expand( $B, V$ )
9:   until  $\max_{\alpha \in V} \alpha \cdot b_0 \geq V_{target}$ 
10:  Return  $V$ 
11: end function

12: function Backup ( $B, V'$ )
13:   $V \leftarrow \{\}$ 
14:  for  $a^R \in \mathcal{A}^R$  do
15:    for  $\alpha'_i \in V'$  do
16:      for  $s \in S$  do
17:         $Q_H(s, a^H) = \sum_{s'} T(s, a^H, a^R, s')$ 
18:           $\alpha_{v(a^H)}(s')$ 
19:      end for
20:       $\Gamma^{a^R} \leftarrow \alpha_i(s) = r(s) + \gamma \cdot$ 
21:         $\sum_{a^H} \pi_H(a^H | Q_H(s, a^H)) \cdot$ 
22:         $\sum_{s'} P(s', a^H | s, a^R, \alpha'_i) \cdot \alpha_i(s')$ 
23:    end for
24:  end for
25:  for  $b \in B$  do
26:     $V_b \leftarrow \{\}$ 
27:    for  $a^R \in \mathcal{A}^R$  do
28:       $V_b \leftarrow \operatorname{argmax}_{\alpha \in \Gamma^{a^R}} \alpha \cdot b$ 
29:    end for
30:     $V \leftarrow \operatorname{argmax}_{\alpha \in \Gamma_b} \alpha \cdot b$ 
31:  end for
32:  Return  $V$ 
33: end function

34: function Expand ( $B', V'$ )
35:   $B \leftarrow B'$ 
36:  for  $b, \alpha \in B', V'$  do
37:     $B_b \leftarrow \{\}$ 
38:    for  $a^R \in \mathcal{A}^R$  do
39:       $s \sim b(s)$ 
40:       $a^H \sim P(a^H | s, a^R, \alpha)$ 
41:       $b'(s') = \eta \sum_s P(s', a^H | s, a^R, \alpha) b(s)$ 
42:      where  $\eta$  is the normalizing constant
43:       $B_b \leftarrow B_b \cup b'$ 
44:    end for
45:     $B \leftarrow B \cup \operatorname{argmax} \|b - b'\|_1, \forall b \in B_b, b' \in B'$ 
46:  end for
47:  Return  $B$ 
48: end function

```

Algorithm 3 Adapted POMCP for CIRL games

```

1: function Search ( $h$ )
2:   repeat
3:     if  $h = \text{empty}$  then
4:        $s \sim I$ 
5:     else
6:        $s \sim B(h)$ 
7:     end if
8:     SIMULATE( $s, h, 0$ )
9:   until TIMEOUT()
10:  Return  $\text{argmax}_{a^R} V(ha^R)$ 
11: end function

12: function Rollout ( $s, h, \text{depth}$ )
13:  if  $\gamma^{\text{depth}} < \epsilon$  then
14:    Return 0
15:  end if
16:   $a^R, a^H \sim \text{Uniform}(A^R), \text{Uniform}(A^H)$ 
17:   $s' \sim T(s, a^R, a^H)$ 
18:  Return  $r(s) + \gamma \cdot \text{ROLLOUT}(s', ha^R a^H, \text{depth} + 1)$ 
19: end function

20: function Simulate ( $s, h, \text{depth}$ )
21:  if  $\gamma^{\text{depth}} < \epsilon$  then
22:    Return 0
23:  end if
24:  if  $h \notin T$  then
25:    Return ROLLOUT( $s, h, \text{depth}$ )
26:  end if
27:   $a^R \leftarrow \text{argmax}_{a^R} V(ha^R) + c\sqrt{\frac{\log N(h)}{N(ha^R)}}$ 
28:   $\theta \leftarrow s_\theta$ 
29:   $a^H \leftarrow \text{SAMPLEHUMANACTION}(\theta, h, a^R)$ 
30:   $s' \sim T(s, a^R, a^H)$ 
31:   $R \leftarrow r(s) + \gamma \cdot \text{SIMULATE}(s', ha^R a^H, \text{depth} + 1)$ 
32:   $B(h) \leftarrow B(h) \cup \{s\}$ 
33:   $N(h) \leftarrow N(h) + 1$ 
34:   $N(ha^R) \leftarrow N(ha^R) + 1$ 
35:   $N(ha^R a^H) \leftarrow N(ha^R a^H) + 1$ 
36:   $V(ha^R) \leftarrow V(ha^R) + \frac{R - V(ha^R)}{N(ha^R)}$ 
37:   $V_\theta(ha^R a^H) \leftarrow V_\theta(ha^R a^H) + \frac{R - V_\theta(ha^R a^H)}{N_\theta(ha^R a^H)}$ 
38:  Return  $R$ 
39: end function

40: function SampleHumanAction ( $\theta, h, a^R$ )
41:   $a^H \sim \pi_H \left( a^H \mid V_\theta(ha^R a^H) + c\sqrt{\frac{\log N_\theta(ha^R a^H)}{N_\theta(ha^R a^H)}} \right)$ 
42:  Return  $a^H$ 
43: end function

```

C. Additional Experiments Omitted in Section 5

C.1. PBVI Experiment on Cooking Domain

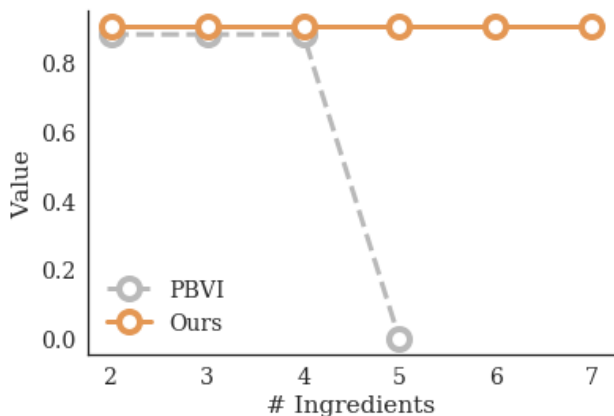


Figure 5. Value attained by PBVI and our approximate algorithm for various numbers of ingredients. (For the first 3 data points, the values attained by both methods were the same – we jittered the plot slightly for visibility).

In our second of the three experiments conducted in Section 5, we compared the values attained by PBVI and our adaptation, with one hour of computation time. We first fixed the number of recipes and varied the number of ingredients. The results of this experiment are presented in Figure 5. We found that both these algorithms, but especially our adapted algorithm, scaled much better with the number of ingredients than their exact VI counterparts. For simpler games, with 3 and 4 ingredients, both algorithms attained the maximal value of 0.9025. However, with 5 ingredients, PBVI found a value of 0 in an hour. In contrast, our algorithm easily solved the game with 5, 6, and 7 ingredients, attaining the maximal value of 0.9025. We next fixed the number of ingredients and varied the number of recipes. Again, our adapted algorithm outperformed PBVI. For example, with 4 recipes, our adapted method attains a value of 0.67875 while the standard PBVI method attains a value of 0.45125.

These results suggest that our modified Bellman update allows PBVI to scale to larger CIRL games, especially in terms of the size of \mathbf{H} 's and \mathbf{R} 's action space. This offers further support to our hypothesis.

C.2. Experiment on *RockSample* Domain

C.2.1. DOMAIN

The second benchmark CIRL domain we present is an extension of the POMDP benchmark domain *RockSample*, that models Mars-rover exploration (Smith & Simmons, 2004). Consider a collaborative task where a human \mathbf{H} wants to take samples of some rocks from Mars, with the help of a rover \mathbf{R} deployed on Mars. There are some number of hours during the day (working hours) when \mathbf{H} can control \mathbf{R} herself but for the rest of the day, \mathbf{R} has to behave autonomously. Not all types of rocks are of equal scientific value; \mathbf{H} knows the value of each of these rocks but \mathbf{R} does not. (Once again, we assume that \mathbf{H} cannot communicate these values to \mathbf{R} directly.)

Formally, consider an instance of *RockSample* on a $m \times m$ grid, with n rocks, each of which belong to one of k different types. The state space is a cross-product of the x- and y-coordinate of \mathbf{R} with n binary features $IsSampled_i = \{Yes, No\}$, which indicate which of the rocks have already been sampled. (Each rock can only be sampled once.)

RockSample is a turn-based game: \mathbf{R} may first take $l_{\mathbf{R}}$ steps in any of the four cardinal direction (during those hours when it is running autonomously) after which \mathbf{H} may similarly take $l_{\mathbf{H}}$ steps (during the remaining hours). Thus, the set of actions available to \mathbf{H} is the set of all trajectories of length exactly $l_{\mathbf{H}}$ while that available to \mathbf{R} is the set of all trajectories with length at most $l_{\mathbf{R}}$. (\mathbf{R} may wait on any specific step if it requires more information while \mathbf{H} may not wait since she has all the information required.)

The set of all reward parameters Θ is composed of a collection of k -dimensional vectors, where the i^{th} entry represents the reward received for sampling rock i . Both agents receive the reward specified by the true reward parameter θ when they sample a rock and receive no reward for any other action.

C.2.2. DETAILS OF EXPERIMENT

We repeat the experiment from Section 5.1 of the main paper on this new domain, with a 5×5 grid ($m = 5$), 3 types of rocks ($k = 4$), and 4 rocks total ($n = 4$).

This domain is much more complex than the cooking domain. For example, note that for even the simplest iteration of this domain, with $l_{\mathbf{H}} = l_{\mathbf{R}} = 1$, $|\mathcal{A}^{\mathbf{H}}| = 4$ and $|\mathcal{A}^{\mathbf{R}} = 5|$ (since \mathbf{R} can also choose to wait); if we raise $l_{\mathbf{R}}$ slightly to 2, we have $|\mathcal{A}^{\mathbf{R}} = 13|$. Hence, we only ran our experiments with POMCP, which scales the best of the three types of the algorithms.

We found similar results to that of Section 5.2 of the main paper. For any value of $l_{\mathbf{R}}$ beyond 1, FV-POMCP and POMCP fail to solve the problem; the branching factor of the search tree they both construct is too large and thus, both methods

deplete the 16GB of memory in our system almost immediately. Our method however manages to scale to larger values of l_H and l_R with relative ease; our method successfully computed the optimal policy for this domain with values of $l_H = l_R = 2$ within two hours of computation.

D. Additional Figure from Section 6.2

Figure 6 describes the results from the experiment in Section 6.2 of the main paper. In this experiment, we analyzed the performance of CIRL and IRL on our collaborative domain when the human does not behave rationally (and when the robot may or may not be aware of this fact).

		Robot's Model of Human				Average	
		CIRL		IRL			
		Rational	Boltzmann-Rational	Rational	Boltzmann-Rational		
Human's Actual Behavior	CIRL	Rational	1	1	1	0.754 ± 0.43	0.9385
	Boltzmann-Rational	0.961 ± 0.16	0.969 ± 0.16	0.972 ± 0.17	0.752 ± 0.44	0.9135	
Human's Actual Behavior	IRL	Rational	0.743 ± 0.46	0.686 ± 0.47	0.698 ± 0.45	0.667 ± 0.47	0.6985
	Boltzmann-Rational	0.357 ± 0.47	0.409 ± 0.49	0.378 ± 0.47	0.484 ± 0.50	0.407	
Average		0.76525	0.766	0.762	0.66425		

Figure 6. The proportion of times that **H** and **R** prepared the correct recipe on the cooking domain when **R** is trained with, and **H** actually behaves according to, a variety of different behaviors. They were significantly more successful at preparing the correct recipe when **H** behaved pedagogically according to CIRL.

E. Results of Follow-up Experiment to Section 6.2

The manipulated variables and dependent measures are exactly as in Section 6.2, with the sole exception being that we considered 10 possible human policies instead of 2. The 10 policies were chosen from a 5×2 factorial of the human's behavior and presence of bias. The 5 possible behaviors were rational, Boltzmann-rational with $\beta = 1$, Boltzmann-rational with $\beta = 5$, ϵ -greedy with $\epsilon = 0.1$, ϵ -greedy with $\epsilon = 0.01$. The 2 possible levels for presence of bias were "No Bias" and "Bias", where "Bias" denoted that **H** had a systematic preference for choosing the "Wait" action. (In our setting, **H** received a reward of 0.25 every time she chose the "Wait" action.)

The results of our experiment are presented below in Figure 7 as a heat map. Much like the simpler experiment in Section 7.2 of the main paper, we find that **H** and **R** are much more successful when **H** behaves pedagogically and that in the presence of pedagogic behavior, the team's performance is more robust to any sub-optimality from **H**.

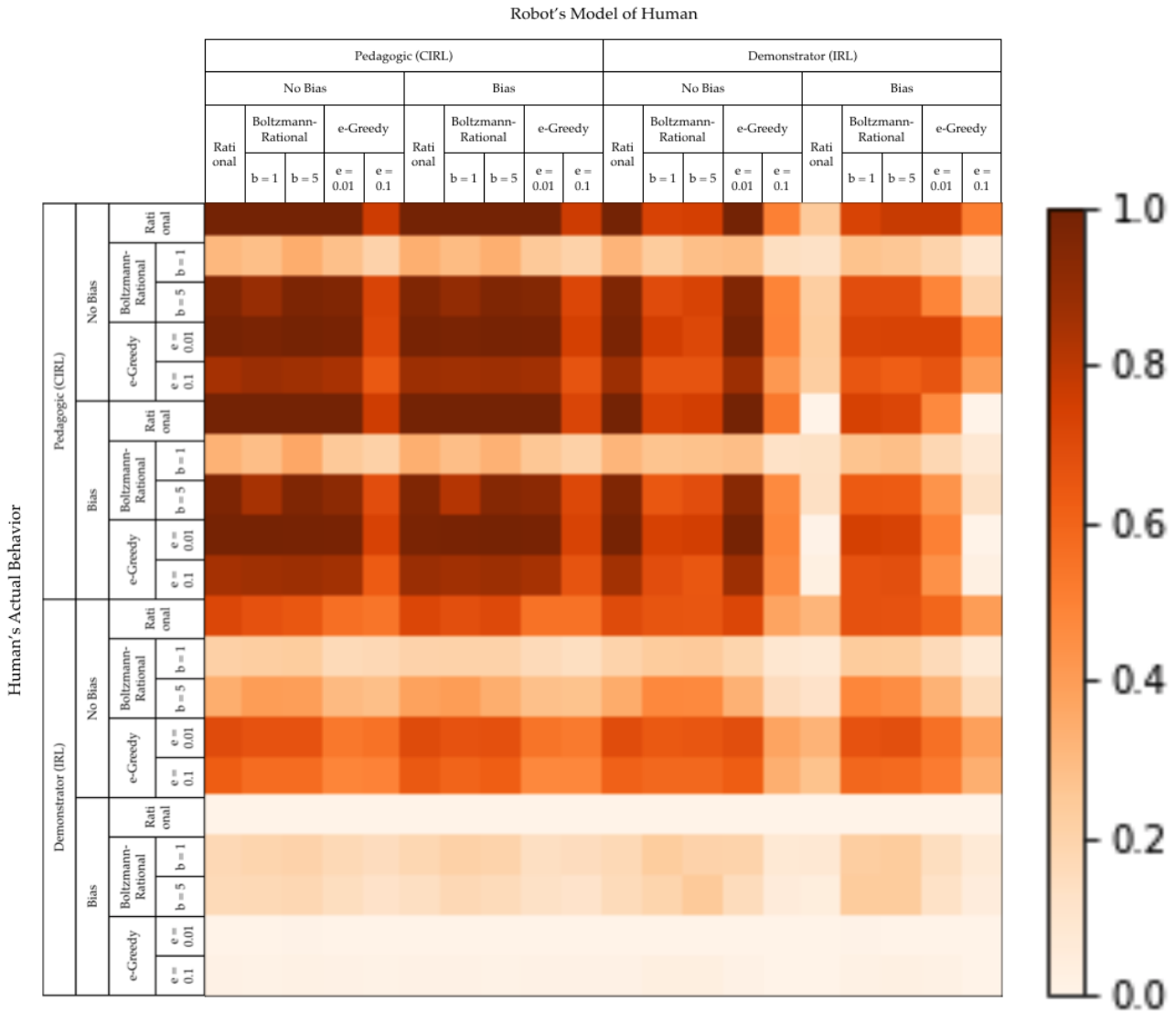


Figure 7. The proportion of times that **H** and **R** prepared the correct recipe on the cooking domain with 3 ingredients and 4 recipes when **R** is trained with, and **H** actually behaves according to, a variety of different behaviors. They were significantly more successful at preparing the correct recipe when **H** behaved pedagogically by following a policy specified by CIRL.

F. Preliminary Human Subjects Study

Previous work observed that CIRL has the potential to outperform standard IRL, and achieve value alignment by allowing a robot to exploit pedagogic behavior from humans (Hadfield-Menell et al., 2016). However, CIRL was only empirically shown to improve upon the performance of IRL in theory or, as in our main paper, in simulation. This does not guarantee that we will observe a similar result in the real world, where the robot interacts with imperfect humans.

Our goal is to investigate whether the benefits of using CIRL over IRL in human-robot collaboration tasks carry over to practice. Here, we conduct a *very preliminary* investigation into whether humans behave pedagogically in practice, and whether a robot trained with CIRL achieves value alignment more successfully than one trained with IRL.

F.1. Hypotheses

We anticipate that humans will objectively succeed at a collaborative task more frequently when collaborating with a CIRL robot as opposed to an IRL robot, especially when the task is complex. We also expect that humans will subjectively prefer to work with a CIRL robot instead of an IRL robot.

H1 - Objective Performance. *The type of algorithm used will positively affect the collaboration objectively across a range of problem difficulties, with CIRL being better than IRL.*

H2 - Objective Performance in Complex Problems. *On more complex problems, the type of algorithm used will positively affect the collaboration objectively, with CIRL being better than IRL.*

H3 - Perceptions of the Collaboration. *The type of algorithm used will positively affect the participants' perception of the collaboration, with CIRL being better than IRL.*

F.2. Experimental Design

To explore the effect of the type of robot on human-robot collaboration, we conducted a counterbalanced within subjects study.

F.2.1. EXPERIMENTAL DOMAIN

Participants collaborated with a virtual robot on the cooking task illustrated briefly in Figure 8 and described extensively in section 5.1 of the main paper. For this experiment, we kept the number of ingredients in the domain fixed at 3, and the length of the horizon fixed at 2.

The robot moved first in this domain. The human was allowed to observe the robot's move before selecting her own move.

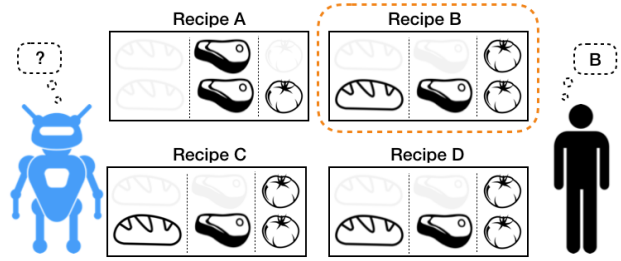


Figure 8. We conduct a *very preliminary* investigation into whether humans do indeed behave pedagogically in practice and whether, as a result, CIRL is more effective than IRL for practically achieving value alignment. Participants collaborated with two robots, one trained with CIRL and another with IRL, to prepare a specified recipe, selected from a larger set. Both the participant and robot were allowed to make a single ingredient at each step but were only given two steps to complete the recipe; so, the human could not succeed without the robot's help. The robot did not know which recipe the participant was instructed to prepare. Participants therefore had to simultaneously teach the robot about their preferred recipe and make progress toward successfully preparing the recipe.

F.2.2. MANIPULATED VARIABLES

We manipulated two variables in our experiment. The first was the *type of robot* used; the two levels were CIRL and IRL. We henceforth refer to a robot trained with CIRL as a CIRL-robot and similarly, to a robot trained with IRL as an IRL-robot.

We additionally wanted to investigate how both robots behaved across a variety of problem difficulties. We suspected that both robots would behave similarly on simple problems due to the straightforward nature of the tasks but that they would behave differently on more complex tasks where they could achieve the goal in a variety of ways. Hence, we additionally manipulated the *number of recipes* used in the task from 2 to 5.

We initially attempted to also vary the length of the horizon on the collaborative tasks. However, we could only solve the longer horizon methods with POMCP, an approximate solver. Hence, there was no guarantee that the solution computed for these problems would be optimal or would be of similar quality across various runs. To avoid confounding the results, we chose to not vary the length of the horizon, and kept it fixed at 2.

We ran a full (2 by 4) factorial experiment with these two manipulated variables, leading to a total of 8 conditions.

F.2.3. PROCEDURE

Participants entered the lab and were administered a pre-study questionnaire. Next, the experimenter explained the collaborative task and informed participants that they would be working with two different robots during the course of the experiment. The experiment was administered virtually

– the participants did not interact with physical robots.

They performed the task four times (one for each possible number of recipes) with one robot chosen at random, and then were administered a questionnaire and asked to describe the robot they had just worked with. They then performed the task four more times with the other robot and were similarly administered a questionnaire. They were finally administered a post-study questionnaire.

F.2.4. PARTICIPANT ASSIGNMENT METHOD

A total of 12 participants (10 males, 2 females, aged 18-25) were recruited from the local community. Ten of the participants reported having a technical background.

The experiment used a within-subjects design because it enables participants to compare the two robots. They were informed that one of the robots was a "student" robot that expected to be taught, and that the other was an "observer" robot that did not expect to be taught but would learn by watching the human perform the task as best as they could. They were made aware of which robot was the "student" and which was the "observer", so that they may behave accordingly and maximize their chance of succeeding at the task.

The order of the robot was counterbalanced to control for order effects. The recipe that participants were instructed to prepare in each condition was randomly chosen from the set of possible recipes to eliminate any systematic or familiarity bias.

F.2.5. DEPENDENT MEASURES

The measures capture the success of a collaboration in both objective and subjective ways, and are based on Hoffman's metrics for fluency in human-robot collaborations (Hoffman, 2013).

The objective measure was *success at preparing the desired recipe*. Participants were assigned a score of one when they succeeded and a score of zero when they failed.

Table 1 shows the six subjective scales that were used, together with a few forced-choice questions. We did not include the questions on Safety/Comfort since the participants did not interact with physical robots. The scales on Robot Contribution and Trust were shortened to avoid asking participants too many questions. The scale on Predictability was rephrased to more appropriately describe the setup of the experiment.

Additionally, participants answered forced-choice questions at the end, about which robot was easier to work with and which robot they preferred.

Fluency

1. The human-robot team worked fluently together.
2. The robot contributed to the fluency of the team interaction.

Robot Contribution [shortened]

1. I had to carry the weight to make the human-robot team better.
2. The robot contributed equally to the team performance.
3. The robots performance was an important contribution to the success of the team.

Trust [shortened]

1. I trusted the robot to do the right thing at the right time.
2. The robot was trustworthy.

Capability

1. I am confident in the robots ability to help me.
2. The robot is intelligent.

Predictability [rephrased for clarity]

1. The robots ingredient selection matched what I would have expected.
2. The robots ingredient selection was surprising.

Forced-Choice Questions

1. Which robot was the easiest to work with?
 2. Which robot do you prefer?
-

Table 2. Subjective Measures

F.3. Analysis

F.3.1. H1 - OBJECTIVE PERFORMANCE

A repeated measures ANOVA on success at preparing the desired recipe showed that CIRL was only marginally better than IRL, when measured across all numbers of recipes ($F(1,11) = 3.667, p = 0.08$). This offers some evidence in support of **H1**.

This is in line with the left plot in Figure 9, which show the results of the experiment. We see that the CIRL-robot outperforms the IRL-robot when averaged across all number of recipes but the improvement is marginal; the error bars for both algorithms have significant overlap.

F.3.2. H2 - OBJECTIVE PERFORMANCE IN COMPLEX PROBLEMS

A repeated measures ANOVA on success at preparing the desired recipe showed that there was a statistically significant interaction effect between the algorithm used and the number of recipes ($F(1,11) = 16.18, p = 0.002$). A post-hoc analysis with Tukey HSD revealed that on complex problems with 5 recipes, CIRL significantly outperformed IRL, but on simple problem, there was no difference in performance between the two algorithms. This offers strong evidence for **H2**.

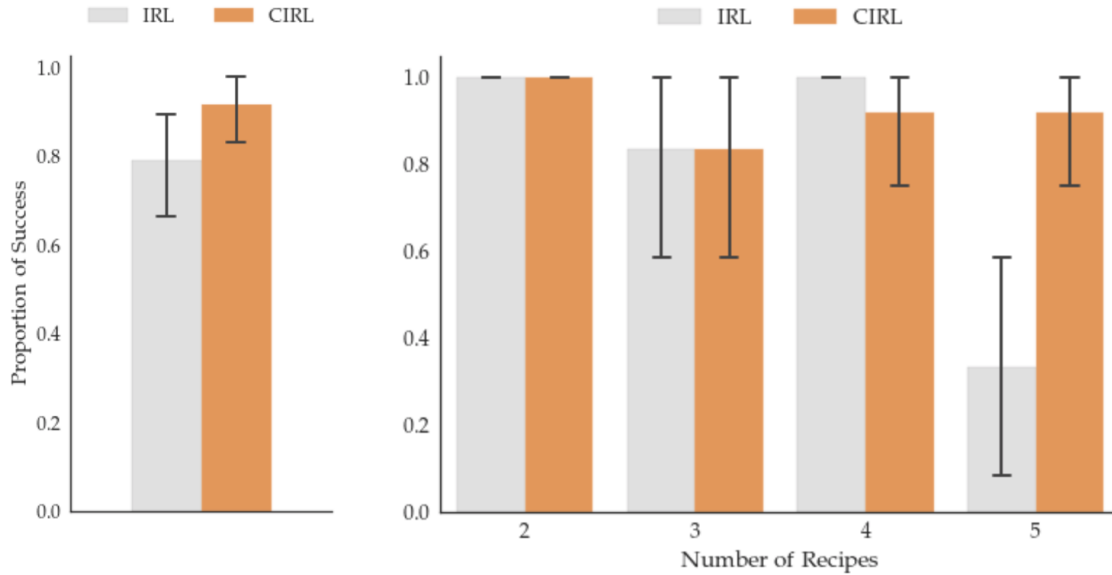


Figure 9. The proportion of all trials on which the human-robot team successfully prepared the correct recipe. (Left) Averaged across all number of recipes: the CIRL-robot only marginally outperforms the IRL-robot. (Right) For each number of recipes: the CIRL-robot only outperforms the IRL-robot for the most complex experiment (with 5 recipes). In all other experiments, both robots perform similarly.

The right plot on Figure 9 echoes these findings. On problems with 2, 3, or 4 recipes, the proportion of trials on which the CIRL-robot and IRL-robot prepared the correct recipe is very similar. However, on problems with 5 recipes, the CIRL-robot was much more successful than the IRL-robot in the collaboration task.

F.3.3. H3 - PERCEPTIONS OF THE COLLABORATION

Scale	Cronbach's α	$F(1,11)$	p-value
Fluency	0.84	23.14	<0.001
Robot Contribution	0.69	17.73	0.001
Trust	0.89	16.95	0.002
Capability	0.81	20.07	<0.001
Predictability	0.86	5.189	0.04
Forced-Choice	0.87	6.494	0.03

Table 3. Results of ANOVA on subjective metrics collected from a 7-point Likert-scale survey.

Table 3 shows the results of the experiment. The internal consistency of each scale is reported via Cronbach's α . All scales except one had "good" consistency, the exception being robot contribution, whose consistency was "questionable" with a Cronbach α of 0.69. Scale items were combined into a score and analyzed with repeated-measures ANOVAs. Figure 10 plots the results.

The score produced by the overall forced-choice questions was significantly affected by the type of robot. The CIRL-robot had a significantly higher score than the IRL-robot

($p = 0.03$); it was rated as being easier to work with by 9 of the 12 participants, and was preferred by 10 of the 12 participants. One participant rated the IRL robot as being easier to work with but preferred to work with the CIRL robot, remarking that he felt more "understood" by the CIRL robot.

All the Likert ratings showed a significant effect for type of robot as well; the CIRL-robot was rated significantly higher than the IRL-robot in *every case* (with $p < 0.01$ in all but one case – predictability). The biggest difference between the two types of robot were in *fluency* and *capability* (both $p < 0.001$). Several participants described the IRL robot as "not intelligent", with one remarking that she felt "the only reason we succeeded as much as we did was because some of the problems were so simple."

These results offer strong evidence in favor of H3.

F.4. Discussion

We have empirically provided strong evidence that suggests that, in practice, CIRL is a more effective framework than IRL for value alignment. In our experiments, participants were objectively more successful at performing the specified human-robot collaboration task when working with a CIRL-robot than with an IRL-robot. Our results further suggest that CIRL leads to more fluent interaction between human and robot; our participants broadly preferred working with the CIRL robot than the IRL robot.

Interestingly, when asked to describe their behavior, many

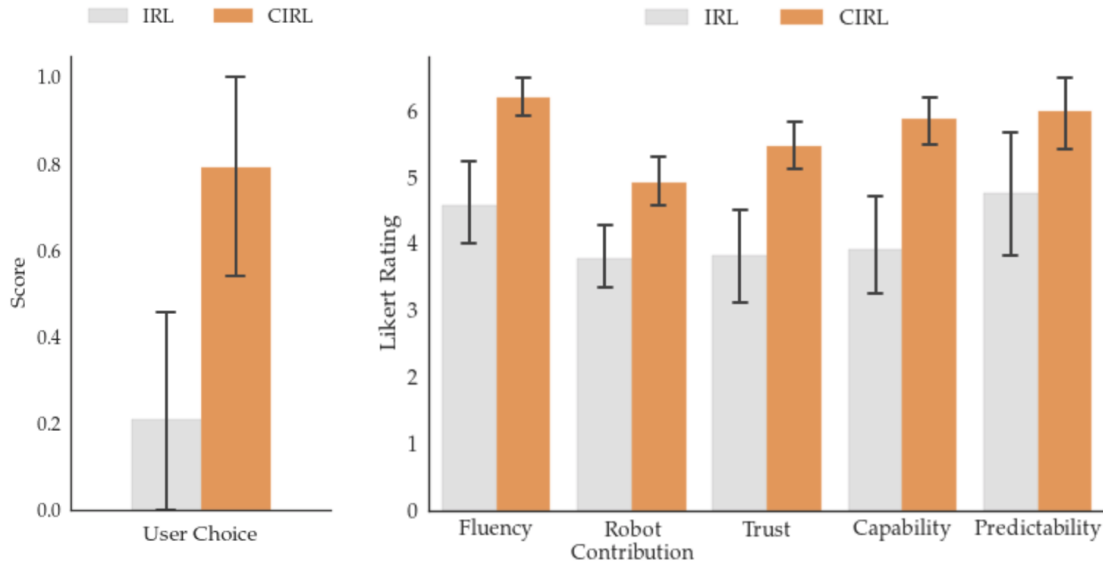


Figure 10. Findings for subjective measures.

participants described behaving similarly with both robots. One participant remarked that regardless of the robot she interacted with, she was "picking her ingredients to eliminate the wrong recipes as quickly as possible."

These remarks agree with the notion that humans tend to behave pedagogically when working with a learner in practice. It is then perhaps no surprise that CIRL significantly outperformed IRL – by exploiting the pedagogic nature of humans, the CIRL-robot was able to infer more information more quickly than the IRL-robot was.

F.4.1. LIMITATIONS AND FUTURE WORK

Due to the computational challenges of CIRL (outlined in the main paper), our experimental domain was still relatively straightforward. The short horizon nature of our task may have made it easier for participants to behave pedagogically and questions remain as to whether people will behave similarly on more complex problems.

Additionally, the demographics of the participants of our survey were rather skewed toward males from technical backgrounds. It is entirely possible that people from technical backgrounds would be more informed about the behavior of the two robots and therefore able to more successfully collaborate with the robots than a non-technical person would.

In future work, we will explore how people behave in collaborative tasks over a long horizon, where their desire or ability to behave pedagogically may be impeded. Furthermore, we intend to deploy our algorithms on real robots and investigate how humans behave in collaborative tasks with actual robots as opposed to virtual ones on computer

screens. To do so, we intend to develop a better online solution method for CIRL with stronger theoretical guarantees on the quality of the solution, thereby allowing us to solve larger problems in practice with real individuals.

Acknowledgements

This work was supported in part by grants from the NSF NRI, and the Open Philanthropy Project.

References

- Amato, C., Dibangoye, J. S., and Zilberstein, S. Incremental Policy Generation for Finite-Horizon DEC-POMDPs. In Gerevini, A., Howe, A. E., Cesta, A., and Refanidis, I. (eds.), *ICAPS*. AAAI, 2009. ISBN 978-1-57735-406-2. URL <http://dblp.uni-trier.de/db/conf/aips/icaps2009.html#AmatoDZ09>.
- Amodei, D. and Clark, J. Faulty Reward Functions in the Wild. <https://blog.openai.com/faulty-reward-functions/>, 2016.
- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., and Mané, D. Concrete Problems in AI Safety. *CoRR*, abs/1606.06565, 2016. URL <http://arxiv.org/abs/1606.06565>.
- Bernstein, D. S., Givan, R., Immerman, N., and Zilberstein, S. The Complexity of Decentralized Control of Markov Decision Processes. *Mathematics of Operations Research*, 27(4):819–840, 2002.

- Bostrom, N. *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press, Oxford, UK, 1st edition, 2014. ISBN 0199678111, 9780199678112.
- Fisac, J. F., Gates, M. A., Hamrick, J. B., Liu, C., Hadfield-Menell, D., Palaniappan, M., Malik, D., Sastry, S. S., Griffiths, T. L., and Dragan, A. D. Pragmatic-pedagogic value alignment. *International Symposium on Robotics Research*, 2017.
- Hadfield-Menell, D., Russell, S. J., Abbeel, P., and Dragan, A. Cooperative Inverse Reinforcement Learning. In *Advances in neural information processing systems*, pp. 3909–3917, 2016.
- Hansen, E. A. Dynamic Programming for Partially Observable Stochastic Games. In *In Proceedings Of The Nineteenth National Conference On Artificial Intelligence*, pp. 709–715, 2004.
- Hoffman, G. Evaluating fluency in human-robot collaboration. In *International conference on human-robot interaction (HRI), workshop on human robot collaboration*, volume 381, pp. 1–8, 2013.
- Kaelbling, L. P., Littman, M. L., and Cassandra, A. R. Planning and Acting in Partially Observable Stochastic Domains. *Artif. Intell.*, 101(1-2):99–134, May 1998. ISSN 0004-3702. doi: 10.1016/S0004-3702(98)00023-X. URL [http://dx.doi.org/10.1016/S0004-3702\(98\)00023-X](http://dx.doi.org/10.1016/S0004-3702(98)00023-X).
- Kocsis, L. and Szepesvári, C. Bandit Based Monte-Carlo Planning. In *Proceedings of the 17th European Conference on Machine Learning, ECML’06*, pp. 282–293, Berlin, Heidelberg, 2006. Springer-Verlag. ISBN 3-540-45375-X, 978-3-540-45375-8. doi: 10.1007/11871842_29. URL http://dx.doi.org/10.1007/11871842_29.
- Kurniawati, H., Hsu, D., and Lee, W. S. SARSOP: Efficient Point-Based POMDP Planning by Approximating Optimally Reachable Belief Spaces. In *In Proc. Robotics: Science and Systems*, 2008.
- Ng, A. Y. and Russell, S. Algorithms for Inverse Reinforcement Learning. In *in Proc. 17th International Conf. on Machine Learning*. Citeseer, 2000.
- Ong, S. C., Png, S. W., Hsu, D., and Lee, W. S. Planning under uncertainty for robotic tasks with mixed observability. *The International Journal of Robotics Research*, 29(8): 1053–1068, 2010.
- Pineau, J., Gordon, G., and Thrun, S. Point-Based Value Iteration: An Anytime Algorithm for POMDPs. In *IJCAI*, volume 3, pp. 1025–1032, 2003.
- Russell, S. and Norvig, P. A Modern Approach. *Artificial Intelligence. Prentice-Hall, Egnlewood Cliffs*, 25:27, 1995.
- Silver, D. and Veness, J. Monte Carlo Planning in Large POMDPs. In *Advances in neural information processing systems*, pp. 2164–2172, 2010.
- Simon, H. A. Models of man; social and rational. 1957.
- Smith, T. and Simmons, R. Heuristic search value iteration for pomdps. In *Proceedings of the 20th conference on Uncertainty in artificial intelligence*, pp. 520–527. AUA Press, 2004.
- Sondik, E. J. *The Optimal Control of Partially Observable Markov Processes*. PhD thesis, Stanford University, 1971.
- Tversky, A. and Kahneman, D. Judgment under uncertainty: Heuristics and biases. In *Utility, probability, and human decision making*, pp. 141–162. Springer, 1975.
- Ye, N., Somani, A., Hsu, D., and Lee, W. DESPOT: Online POMDP Planning with Regularization. 58:231–266, 2017.