

# 1 Introduction

Much of human social life occurs in contexts where people must coordinate their actions with those of others. From party planning to flying a rocket ship to the moon, groups of people have accomplished great things by reasoning as a team and engaging in jointly intentional behavior [45]. Indeed, some have argued that, because most other animals lack the capacity to work adaptively as cohesive unit across many domains, team reasoning may be the hallmark of human sociality [50].

Artificial agents are becoming ubiquitous in everyday life. For these agents to collaborate effectively with humans, however, they will need behave in human-like ways. Unfortunately for agent designers, there is no simple behavioral rule (or set of rules) that could be prescribed to cover all contexts, communities, and cultures. People perpetually learn new rules and alter their presumed roles,<sup>1</sup> and so must artificial agents.

Most attempts to build collaborative, artificial agents begin with a model (however flawed) of human behavior. One possibility is to assume we are *Homo economicus*—a term from the 1800s that portrays individuals as being guided by self interest, and taking actions without regard to their impact on the collective. Alternatively, one might assume an individual to be *Homo sociologicus*—a more recent term depicting individuals as pursuing goals imprinted on them by society.

The homo-sociologicus model was not intended as a complete description of human behavior, which is clearly self interested at times. But, psychologists and behavioral economists have also more or less rejected the homo-economicus model [?, ?], proposing countless homo-sociologicus-leaning models in its stead [?]. Consistent with these findings, this proposal is founded on the assumption that **social preferences** hold the key to successful collaborations. Social preferences are in part **distributive**, meaning they evaluate how resources are allocated, e.g., efficiently or equitably; and they are in part **reciprocal**, meaning they evaluate the perceived intentions of the actors in the environment.

We plan to model interactions among such social agents, both human and artificial, in the framework of game theory. Following Bacharach [?], not only do we assume that preferences are social rather than individual, we go one step further and assume that decision making takes the form of “What should *we* do?” rather than “What should *I* do?” In other words, we move from the usual assumption of best-reply reasoning [12] in games to one of **team reasoning** [?].

We call optimal decision making, when agents hold social preferences and employ team reasoning, **socially rational** behavior. Social rationality is a blend of the *Homo economicus* and *Homo sociologicus* concepts: agents optimize a **social utility function** (i.e., a representation of social preferences), which is sufficiently rich to incorporate perceived (possibly inferred, possibly imagined) societal benefits.

We contend that the computational process by which collaborative behavior emerges strongly depends on social rationality, coupled with learning via repeated interaction. Specifically, the computational process is thus: 1. agents represent social preferences using a social utility function tailored to the environment and the actors;<sup>2</sup> 2. based on their estimate of this function, agents use team reasoning to make plans to act; 3. agents act, simultaneously or sequentially, as is appropriate in their environment; 4. agents update their estimates of their social utility functions. Steps 2 through 4 repeat so long as the interaction continues.

When agents make a plan that optimizes a social utility function, they are choosing actions that jointly optimize the behavior of the collective. When agents update their estimate of the social utility function, they do so using the observed history of joint actions in the environment. If this process stabilizes, the social utility function represents the preferences of the community as a whole, and the equilibrium implements a joint plan of action that optimizes that social utility function. Indeed, in our computational model, emergent collaborative behavior and social preferences are expected to reinforce one another.

Classical game theory assumes that utility functions are exogeneously determined: i.e., agents come to the table with fixed extrinsically-set utilities. However, there is broad agreement in the judgment and decision-making community that preferences are *constructed*, not elicited [53], in response to available choices. Likewise, our computational framework does not assume that social preferences exist extrinsically; rather, they are constructed in response to observed behaviors. In our conception of social rationality, agents *learn* social preferences through repeated play of the game.

---

<sup>1</sup>As Shakespeare said, “one man in his time plays many parts.”

<sup>2</sup>This utility function is described in terms of features of the environment and the actors; as such, it is applicable to other environments and actors with similar features.

A social utility function is an immediate generalization of the usual notion of an individual’s utility function, but it ascribes value to the outcome as it relates to all actors in the environment. As such, it can be used to capture important elements of an artificial agent’s decision making, including (1) how it values the utility of others [29], (2) how it values the utility of intermediate goals [33], and (3) how it values sanctions against others who exhibit anti-social behavior [?].

We further assume that social utilities can be broken down into two components—an objective component, which is usually a direct function of the rules of the game, and a subjective component, which captures notions of distributivity and reciprocity. Often, the rules of a game alone give rise to multiple (objective) equilibria, so that agents face a difficult equilibrium-selection problem (e.g., [44]) in the underlying game. We contend that socially rational behavior, in which agents optimize a learned social utility function that reflects (constructed) social preferences, can provide a solution to this challenging, and often elusive, problem.

*Under the assumption that human agents are, perhaps boundedly but nonetheless ideally, socially rational creatures, we propose to design and build socially rational artificial agents that learn through repeated play, with the aim being for such agents to collaborate effectively with humans.*

In this proposal, we set out to demonstrate our potential for success, by showing that our computational model is viable. To this end, we present a preliminary set of simulation results showing how reinforcement-learning agents, interacting repeatedly in a grid-game environment, can represent social preferences using a social utility function, so that eventually their emergent collaborative behavior and their discovered social preferences reinforce one another. We then go one step further and show that our approach can also learn (in an off-line, batch fashion) from human–human trace data so that the collaborative behavior that is learned mimics that which is expressed in the trace data. While these two experiments (on machine–machine, interactive and human–human batch data) are indeed promising, it remains to extend these techniques to human–machine interactions, preferably in experimental worlds with embodied robots. **\*\*Amy: HUMAN TRACE DATA!\*\***

The ultimate goal in developing artificial agents that act intelligently in multi-agent scenarios is to apply them to real-world problems. To a limited extent, this goal has already been achieved: artificial agents trade stocks and bid in online ad auctions. In these two market environments, rationality may well be an appropriate model of behavior for the participating agents. However, an artificial agent that plans its collaboration by assuming the humans in its environment will act selfishly and/or are driven by extrinsic utility functions is unlikely to be successful in its collaborations [9, 24]. Instead of individual rationality as the underlying principle guiding human behavior, this proposal is grounded in the assumption that humans are socially rational beings whose preferences and behaviors emerge and reinforce one another through repeated interactions. For artificial agents to successfully collaborate with humans, they should be thus, too!

## 2 Grid Games

A grid game is a game played by two agents on a grid, in which each agent has a goal. See, for example, Figure 1a, which is a 3x5 grid in which the two agents’ initial positions are one another’s goals: Orange begins in position (1,2), Blue’s goal; and Blue begins in position (5,2), Orange’s goal. We refer to grid positions using  $x$ - $y$  coordinates, with (1,1) as the bottom left position.

One grid game proceeds in rounds, each round consisting of multiple turns. On each turn, the agents choose one of five actions (N, S, E, W, or wait), which are then executed simultaneously. In the most basic setup, agents transition deterministically, and there is no tie-breaking when two agents collide.<sup>3</sup> Instead, if their chosen actions would result in a collision with one another, neither agent moves. A round ends when either (or both) players move into their goal, or when a maximum number of turns has been taken.

We present five two-agent grid games, specifically designed to vary the level of coordination required, while at the same time allowing agents to defend against uncooperative partners. In these games, one approach is for an agent to cooperate blindly with its opponent by simply moving out of the opponent’s way, and hoping the opponent then waits for the agent to catch up. However, such strategies can be exploited by uncooperative ones that proceed directly to the goal as soon as their path is unobstructed.

<sup>3</sup>It is a simple matter to vary these rules within our infrastructure, as future experimental design might dictate.

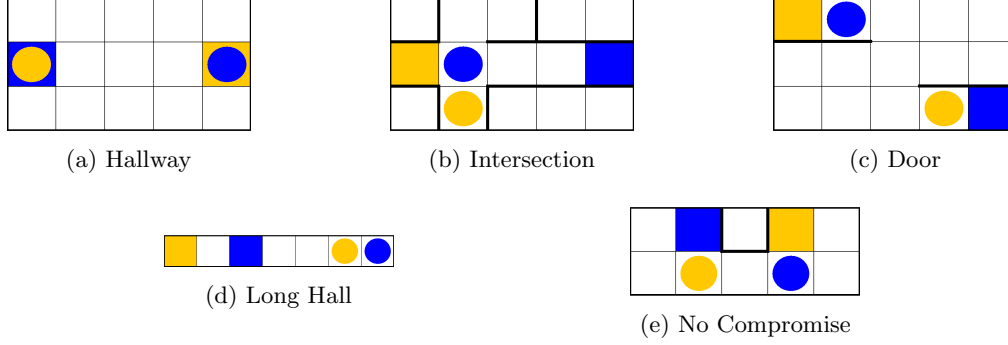


Figure 1: Five example grid games.

To distinguish “unsafe” from “safe” cooperation, we devised a new categorization for strategies in our grid games. We call strategies that allow for cooperation, while at the same time maintain a defensive position in the event that the other agent is uncooperative, **cooperative defensive** strategies. More formally, an agent’s strategy is **cooperative** (C) if it is one that allows both it and its opponent to reach their goals, while an agent’s strategy is **defensive** (D) if its opponent does not have a strategy that allows it to reach its goal strictly first. A cooperative defensive (CD) strategy is both cooperative and defensive.

We now proceed to describe a sample set of grid games, and equilibria comprised of CD strategies (when they exist), to illustrate the kinds of interactions we plan to study. Our first example, Hallway, is depicted in Figure 1a. This game is one in which the agents can choose to coordinate, for example, if both agents agree upon a joint strategy where one agent moves along the top row and the other along the bottom, without interfering with one another. But, an agent could choose to “defect” from this joint strategy, by proceeding straight to its goal. There are CD strategies, however, which defend against such non-cooperative behavior.

For example, if Orange moves south initially to (1,3) and Blue moves west to (4,2), Orange might choose to return and remain on its goal until Blue retreats to (4,3) or (4,1), at which point the players are equidistant from their goals, and both can reach them safely. This joint strategy is an equilibrium comprised of CD strategies, since Orange and Blue both remain in positions where they have the ability to block their opponents until they both have unobstructed equidistant paths to their respective goals.

The grid in Figure 1b (Intersection) requires Blue to defend against the possibility of Orange behaving uncooperatively, which it can achieve by squatting on the orange goal. Orange can then move to (3,1) where both agents are equidistant from their goals. Therefore, this game also has an equilibrium comprised of CD strategies for both players.

This equilibrium is not the shortest path, however. Purely cooperative agents in this game could adopt a joint strategy in which Blue moves east, while Orange waits a single step, before both agents proceed into their goals. This strategy profile is not defensive for Blue though, because it does not have the opportunity to observe if Orange will cooperate (wait) or defect (go north), and therefore cannot defend itself if Orange decides to head straight toward its goal.

Figure 1c (Door) is a grid that requires coordination to navigate through the narrow center space at (3,2). Any equilibrium comprised of CD strategies for this grid must be asymmetric, because it requires one agent to cede to the other agent the center cell. For example, if Orange chooses to cede that cell, it should step west into (2,3) while Blue steps south into (3,2). Then, Orange needs to step east back into (3,3) to prevent Blue from marching straight into its goal. Only when Blue agrees to step aside to (2,2) will they both be equidistant from their respective goals and in position to cooperate. This intricate pattern of first giving way to the opponent, and then forcing them to step around later represents an equilibrium comprised of CD strategies, since both agents are able to prohibit their opponent from reaching the goal first, but still leaves open the possibility for them both to reach their goals, cooperatively.

In the grid in Figure 1d (Long hall), Blue begins one step closer to its goal than Orange does. However, Orange can squat on the blue goal until Blue chooses to cooperate by taking one step back. If Orange can predict when Blue steps back, then Orange can take one step closer to its goal while Blue steps further away, in which case only Orange would reach its goal. The strategy that minimizes the risk to either agent requires

that Blue wait one turn initially, while Orange moves toward its goal. These two strategies comprise a CD equilibrium.

Our last grid, shown in Figure 1e (No compromise), requires not only cooperation, but both agents must also exhibit trust for one another, or both agents cannot arrive at their goals at the same time. For example, Orange may sit on Blue’s goal so that Blue can move to (1, 2). Then, Blue must wait two turns before both agents are equidistant from the goals. If Blue defects and moves south into (1, 1) while Orange moves south into (2, 2), Orange still has the opportunity to go back up north to block Blue from reaching its goal. However, if Blue moves south into (1, 1) when Orange steps east into (3, 2), Blue will arrive at its goal sooner. Therefore, a trust spanning multiple rounds is required for the agents to effectively cooperate in this game.

No equilibrium in CD strategies exists for No Compromise. The game is like Door in that only one player can go move through the middle cell at a time. Unlike Door, however, it is not possible for the agents to simultaneously maintain a defensive position and to signal cooperation, because any cooperative move leads to an asymmetric situation in which the agents are no longer equidistant from their goals. As a result, after one agent cooperates, there is always an incentive for the other agent to defect, and there is nothing the cooperative agent can do to defend itself. Note however, that if Orange sits on the blue goal while Blue walks to (1, 1) and then Blue cooperates by waiting for Orange to walk to (1, 3), Blue’s policy is CD. Still, Orange cannot respond in kind with a CD strategy; the aforementioned strategy is C.

Taking these five grid games as an initial testbed, we performed two pilot studies: the first involved simulations of artificial agents playing against one another; the second pitted humans against other humans on Mechanical Turk [\[1\]](#). **\*\*Amy: add citation to workshop paper\*\*** We summarize the results of these preliminary studies next.

**Simulation Experiments** We carried out a set of simulation experiments with  $Q$ -learning in the grid games presented. We denote the outcome of a round using a pair of letters, where **G** means the agent reached the goal and **N** means the agent did not reach the goal. The first letter in the pair represents the agent’s own outcome and the second represents its opponent’s outcome. For example, **GN** is used to denote that the agent reaches its goal while its opponent does not.

As two  $Q$ -learning algorithms are not guaranteed to converge in self-play, we arbitrarily stopped the learning after 30,000 rounds, and checked the strategies learned. In spite of  $Q$ -learning not explicitly seeking outcomes with high social welfare, it very reliably identified cooperative strategies leading to **GG** outcomes.

In a second set of  $Q$ -learning experiments, we examined the impact of endowing agents with **social preferences**. That is, their utilities no longer depend solely on their own successes and failures—they depend on those of other agents in the environment as well.

To make this idea more precise, we draw a distinction between **objective** and **subjective** reward.<sup>4</sup> Objective reward is the usual reward signal provided to the agent from the environment, and by which behavior is judged. Standard reinforcement-learning agents, such as  $Q$ -learners, seek to optimize their own objective reward signal—we call this preference the **selfish** preference because these agents are only concerned with their own outcomes.

In some environments, however, it is useful to distinguish this objective reward from subjective reward—the perceived quantity the agent is optimizing. Previous work has shown that optimizing subjective reward can sometimes lead agents to be more effective in their acquisition of objective reward [47]. (Indeed, we reach this same conclusion in our experiments.)

Considering the four different outcomes in these games—**GG**, **GN**, **NG**, **NN**—there are 75 different possible preference orderings (allowing for ties). The selfish ordering that ignores the opponent’s outcome is one of these: **GG**  $\sim$  **GN**  $\succ$  **NG**  $\sim$  **NN**. Nine of the 75 orderings are consistent with the selfish ordering, strictly preferring **Gx** to **Nx** for all **x**.

Of particular interest is the **fair** preference, which we define as the objective reward of the agent minus 25% of the difference between its own and the opponent’s objective rewards:  $r_s = r_a - 0.25|r_a - r_o|$ , where  $r_s$  is the agent’s subjective reward,  $r_a$  is the agent’s objective reward, and  $r_o$  is the opponent’s objective reward.<sup>5</sup> By incorporating this fairness term into the agent’s subjective reward, it strictly prefers the following ordering: **GG**  $\succ$  **GN**  $\succ$  **NN**  $\succ$  **NG**. That is, the agent prefers making it to its goal as opposed to

<sup>4</sup>We use the terms “utility” and “reward” interchangeably.

<sup>5</sup>Other percentages would achieve the same result.

not, but it additionally prefers that the opponent only get to its goal if the agent itself does as well. To say it another way, a fair agent wants its opponent to win with it or lose with it.

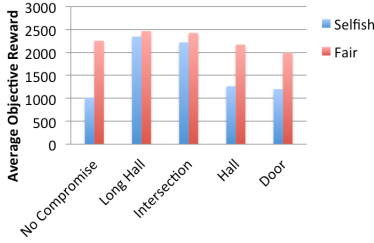


Figure 2: Average score of  $Q$ -learners in self play (30,000 rounds)

Figure ?? shows the result of the selfish and fair agents playing against others of the same type in each of our test grid games. Of the nine orderings, only three (all variations of the fair preference ordering) achieve consistent cooperation in self play across all five grid games. Consequently, fair agents obtain higher total (objective) rewards than others. We also ran all nine preference orderings against one another. The average scores (across both players) in games involving a fair agent tend to be higher than the average scores in games not involving a fair agent. We also analyzed the types of strategies learned by fair and selfish  $Q$ -learners after multiple simulations of various configurations. Interestingly, we found that  $Q$ -learners with fair preferences tend to find CD strategies more often, especially when paired with selfish agents.

In summary,  $Q$ -learners naturally learn to cooperate in the grid games studied, discovering equilibria comprised of CD strategies when they exist. Moreover, cooperation can be induced in these games by manipulating the  $Q$ -learners’ subjective rewards to incorporate the success of others.

**Turk Experiments** We also conducted analogous experiments with humans playing grid games on Mechanical Turk. The results involving human subjects playing grid games are consistent with the results obtained in our simulations of  $Q$ -learners playing grid games. In both sets of experiments, more cooperation was achieved when the treatment incorporated social preferences.

While the subjective utility functions of artificial agents are within our control, so that we can perhaps lead artificial agents towards collaborative behavior, the subjective utility functions of humans are not. Nonetheless, behavioral economists often infer approximations of utility functions from experimental data by assuming some underlying structural form, and then estimating the relevant parameters [7, 16]. **\*\*Amy: CHECK REFERENCES!\*\*** Likewise, one of the primary intended uses of our Turk experimental framework is to generate trace data of humans playing grid games and learning collaborative behavior (such as trust, CD, etc.), so that we can then proceed to infer utility functions and relate them back to the behaviors they produce.

### 3 Computational Framework

Recall, from the Introduction, our computational model: 1. agents represent social preferences using a social utility function; 2. based on their estimate of this function, agents use team reasoning to make plans to act; 3. agents act, simultaneously or sequentially, as is appropriate in their environment; 4. agents update their estimates of their social utility functions.

While certain aspects of social preferences (e.g., fairness) are broadly applicable, Step 1 generally requires a representation that is tailored to the specific environment and actors within. Likewise, Step 3 involves a straightforward simulation in grid games, but is otherwise specific to the robots and their environment.

Only Steps 2 and 4 are sufficiently generic to apply across environments. Step 2 is a **planning** step; as planning is a well-defined optimization problem, it can be accomplished using any off-the-shelf planner for stochastic sequential decision making environments [4, 5, 8, 11, 25, 26]. Step 4, on the other hand, is not a solved problem; indeed, in this proposal, we describe an innovative approach to this problem.

Step 4 can be thought of as a **learning from demonstration** problem [1], in which an agent learns how to behave in an environment by observing an expert. That is, an agent is presented with a data set consisting of multiple examples of behaviors (e.g., state-action trajectories through a Markov decision process), and is tasked with the objective of learning a policy capable of (closely) reproducing the data.

Two widely studied approaches to this problem include **imitation** and **intention** learning. Imitation learning is a supervised learning approach in which an agent learns a mapping from states to actions, with the expert’s behavior serving as training data [39]. After training, a learner follows its learned policy directly.

Learning intentions [33], in contrast, is often framed as an **inverse reinforcement learning** (IRL) problem [2, 38]. Here, the agent’s goal is to learn rewards that motivate the expert’s behavior, after which it applies a planning algorithm to derive a policy that is consistent with those learned rewards. A strength of intention learning is that even simple reward functions can capture complex behavior, leading to greater generalization capabilities (i.e., appropriate behavior in as-yet-unforeseen states) than the more direct approach taken by imitation learning.

We are proposing to use IRL in Step 4. Preliminary experiments, described herein, with one off-the-shelf IRL algorithm, demonstrate the potential of our computational model to facilitate collaborative behavior among artificial agents. But to better support our goal of human–machine collaboration, we also propose a new approach to IRL in which demonstrations are not, by default, assumed to have been generated by an expert. With this additional power, our IRL approach can be applied to all past interactions among agents, separating the successful from the unsuccessful. Indeed, in the computational process we put forth, learning is online, not off, so not all demonstrations should be interpreted as successful. On the other hand, as past interactions are shared context across agents, ideally they would *all* provide crucial guidance for identifying cooperative behavior, and that is what we set out to achieve.

### 3.1 Mathematical Models

As grid games are examples of stochastic games [15, 46], we use this game-theoretic framework to model agent interactions. Like econometricians (e.g., [43]), we assume a structural form for utility functions; and because we are interested in building socially rational artificial agents, we endow our agents with a (potentially) social utility function. We then develop a batch learning algorithm by which an artificial agent can, by learning from demonstrations, recover a social utility function that supports exactly one of many equilibria (assuming the demonstrations exhibit exactly one of many equilibria). We also develop an interactive learning algorithm by which two agents can simultaneously learn a social utility function, and converge to one among a plethora of equilibria, thereby endogenously solving the equilibrium selection process.

**Markov decision processes** A **Markov Decision Process** (MDP) is a model of a single-agent decision making problem, defined by the tuple  $(S, A, T, R)$ , where  $S$  is the set of states in the world;  $A$  is the set of actions that the agent can take;  $T(s' | s, a)$  defines the transition dynamics: the probability the environment transitions to state  $s' \in S$  after the agent takes action  $a \in A$  in state  $s \in S$ ; and  $R(s, a, s')$  is the reward function, which returns the reward the agent receives when environment transitions to state  $s'$  after the agent takes action  $a$  in state  $s$ .

The goal of planning or learning in an MDP is to find a policy  $\pi : S \rightarrow A$  (a mapping from states to actions) that maximizes the expected future discounted reward under that policy:  $E^\pi [\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1})]$ , where  $\gamma \in [0, 1]$  is a discount factor specifying how much immediate rewards are favored over distant rewards.

To find an optimal policy, many algorithms compute the optimal state ( $V^*(s)$ ) and state-action ( $Q^*(s, a)$ ) value functions that specify the expected future discount reward under the optimal policy from each state, and from taking an action in a state and then following the optimal policy respectively. Given these functions, the optimal policy is derived by taking an action in each state with the maximum Q-value:  $\pi(s) \in \arg \max_{a \in A} Q(s, a)$  [6].

**Inverse reinforcement learning** Although there are multiple IRL formalizations and approaches, they all take as input an MDP together with a **family** of a reward functions  $R_\Theta$  that is defined by some parameter space  $\Theta$ , and a dataset  $D$  of trajectories (where a trajectory is a finite sequence of state-action pairs:  $\langle (s_1, a_1), \dots, (s_n, a_n) \rangle$ ). The algorithms then output a specific parameterized reward function  $R_\theta \in R_\Theta$ , which induces a policy that is consistent with the input trajectories. Different IRL algorithms frame the objective function for policy consistency differently. One common approach is to treat the search problem as a probabilistic inference problem [2, 31, 41, 57].

For example, in the maximum-likelihood setting [2], the goal is to find a reward function parameterization that maximizes the likelihood of the data:

$$\theta \in \arg \max_{\theta} L(D | R_\theta) = \prod_{t \in D} \prod_i^{|t|} \pi_\theta(s_i, a_i), \quad (1)$$

where  $\pi_\theta(s, a)$  is a stochastic policy defining the probability of taking action  $a$  in state  $s$  when the parameterized reward function to be maximized is  $R_\theta$ . Typically, the Boltzmann (softmax) stochastic policy over the  $Q$ -values is used:  $\pi_\theta(s, a) = \frac{e^{\beta Q_\theta(s, a)}}{\sum_{a' \in A} e^{\beta Q_\theta(s, a')}}$ , where  $Q_\theta(s, a)$  is the  $Q$ -function when the reward function is parameterized by  $\theta$ . Popular methods for solving Equation 1 (i.e., for carrying out maximum-likelihood IRL (MLIRL)) include gradient ascent and expectation-maximization.

In Bayesian IRL, the goal is to compute a posterior distribution over reward function parameterizations:  $\Pr(R_\theta \mid D) \propto \Pr(D \mid R_\theta) \Pr(R_\theta)$ . Bayesian IRL algorithms typically rely on Markov chain Monte Carlo methods [40, 41]. A related approach is maximum-entropy IRL [56]. More recently, work has also commenced on nonparametric Bayesian IRL methods [10, 36].

IRL algorithms are typically computationally demanding because they require planning in their inner loops. Compatible with all of the above algorithms, **receding horizon IRL** (RHIRL) [33] addresses this limitation by replacing the usual infinite-horizon policy with a receding horizon controller (RHC) that only plans out to some finite horizon from any given state, thereby bounding computation time. An RHC tends to steer an IRL algorithm towards a **shaping reward function** [37] short-term rewards that guide the agent without it having to plan too far ahead. Consequently, RHIRL allows us to plan with a short horizon, thereby saving on computation time.

Our computational framework, and the generalized model of IRL we introduce in Section 3.3, are all agnostic about the particular approach to IRL taken. We plan to experiment with multiple variations.

**Stochastic games** The stochastic games formalism can be viewed as an extension of MDPs to the multi-agent case [28]. A **stochastic game** is defined by the tuple  $(I, S, A^I, T, R^I)$ , where  $I$  is an index set of agents in the environment;  $S$  is the set of states of the environment;  $A^I$  is set of actions for each of the agents with  $A^i$  denoting the action set for agent  $i \in I$ ;  $T(s' \mid s, j)$  is the transition dynamics specifying the probability of the environment transitioning to state  $s' \in S$  when the *joint action*  $j \in \times_i A^i$  of all agents is taken in state  $s \in S$ ; and  $R^I$  is a set of reward functions for each agent with  $R^i(s, j, s')$  denoting the reward received by agent  $i \in I$  when the environment transitions to state  $s' \in S$  after the agents take joint action  $j \in \times_i A^i$  in state  $s \in S$ .

The goal in a stochastic game is to find a joint strategy that satisfies some solution concept. Different solution concepts for stochastic games have been explored in the past including minimax, Nash, correlated, and CoCo equilibria [19, 23, 27, 58]. There are problems with these approaches, however. First, the resulting planners must solve for game-theoretic equilibria in an inner loop, a problem which in the case of Nash equilibrium, for example, is notorious for its computational intractability [13]. Second, in the general case of non-constant sum games, none of the planners that make reasonable assumptions about agent behavior yield unique joint plans, and none has solved the ensuing equilibrium selection problem suitably.

Immediately generalizing from the case of MDPs, the goal of inverse reinforcement learning in stochastic games is to learn a set of reward functions for the stochastic game that describe an environment that would motivate the agents to behave in a way that is consistent with the observed behavior under some solution concept such as a Nash equilibrium [42]. This problem, however, is exceedingly difficult, because the planning that is necessary in the inner loop of an IRL algorithm is subject to the challenges identified by the equilibrium planners mentioned above. In this project, we plan to develop IRL algorithms that facilitate equilibrium selection in a stochastic game.

**Social Reward Functions** In our preliminary studies, we assume players’ individual (game) rewards are known, and our goal is to recover a **social reward function** that combines these known rewards in such a way as to capture social preferences expressed in the joint play of the agents (when one exists).

Like the players’ reward functions, a social reward function operates on states, joint actions, and next states:  $R^S : S \times \times_i A^i \times S$ . In our initial model, we assume the social reward function can be written as linear combination of a **team function**, which represents team goals, and a family of what we call **bias functions** ( $B_\Theta$ ) defined by some parameter space  $\Theta$ . The team function takes as input a multi-agent reward function ( $R^I$ ), and returns a single numeric “team” value for any state-joint action-next state triple. One example of such a function is total welfare (i.e., the sum of all agent rewards):  $\mathcal{T}(R^I, s, j, s') = \sum_i R^i(s, j, s')$ . The bias function family is similar in nature to a reward function family that would be input to a classic IRL algorithm, but operates on state-joint action-next state triples, thereby encoding a bias that motivates collaborative

---

**Algorithm 1** Batch\_Learning( $(I, S, A^I, T, R^I), \gamma, D, \mathcal{T}, B_\Theta$ )

---

**Require:** stochastic game  $(I, S, A^I, T, R^I)$ ; discount factor  $\gamma$ ; multi-agent demonstrations  $D$ ; team function  $\mathcal{T}$ ; and parameterized bias function family  $B_\Theta$ .

$A^M := \times_i A^I$  ▷ Joint action set  
 $R^M(s, a, s') := \mathcal{T}(R^I, s, a, s')$  ▷ Team reward function  
 $R_\Theta^S(s, a, s') := R^M(s, a, s') + B_\Theta(s, a, s')$  ▷ Social reward function family  
 $R_\theta^S := \text{IRL}((S, A^M, T, R_\Theta^S), \gamma, D)$   
**return**  $R_\theta^S$  ▷ Learned social reward function

---

behavior in games. In addition to the parameters of the bias function, the social reward function may also include a parameter to trade-off the team rewards against the biases.

### 3.2 Collaborative Learning Algorithms

Given our representation of social preferences by a social reward function that is parameterized by a family of bias functions, our approach to collaborative learning is simply to optimize those parameters to best match observed behavior. We consider two forms of collaborative learning. The first is **batch** learning, in which an agent learns offline from batches of demonstrations of other agents playing collaboratively. The second is **interactive** learning in which agents learn the social preferences of the collective via repeated interactions. At the core of both algorithms is team reasoning.

**Team Reasoning** In our model, the demonstrations agents learn from are trajectories of state-joint action pairs, specifying the behavior of all the agents in the environment. Likewise, their objective, as team reasoners, is to produce a plan for what all the agents in the game should do. To produce such a plan, the stochastic game must be transformed into an MDP. For the most part, this transformation is straightforward: the states are the same, the MDP action set is the space of joint actions, and the MDP transition dynamics still operate on joint actions (which is the action set in the MDP). A family of social reward functions is then defined as a linear combination of the input team function (which depends on the stochastic game’s reward functions), and the family of bias functions.<sup>6</sup>

In addition to the human-behavioral motivations for building agents that employ team reasoning, team reasoning also makes the problem tractable, because IRL in MDPs is becoming increasingly easier, while IRL in stochastic games remains enormously challenging.

**Batch Learning** The batch learning setting is similar to the standard IRL problem in that the algorithm is given a batch of demonstrations and a family of reward functions over which to search. In our case, however, the batch data consists of trajectories of states and joint actions in a stochastic game, and the reward function family is a social reward function family. To learn the parameterization of the social reward function family that best captures the behavior exhibited in the demonstrations, we convert the stochastic game into an MDP, and then use IRL to learn a parameterization of the reward function family in this MDP.

Pseudocode for the batch learning algorithm is shown in Algorithm 1. The algorithm takes as input a stochastic game; a discount factor; a set of multi-agent demonstrations (which perhaps adheres to some equilibrium: i.e., collaborative behavior); a team function; and a family of bias functions defined by some parameter space  $\Theta$ . The first step of the batch learning algorithm is to transform the stochastic game into an MDP, and the second to apply IRL to that MDP. After learning a social reward function, any single agent can then behave by selecting a joint action from the policy that results from it, and then selecting their individual action from the joint. That is, agent  $i$  takes action  $j^i \in A^i$ , where  $j \in \arg \max_j Q_{R_\theta^S}(s, j)$ .

**Interactive Learning** In the interactive learning setting, an agent plays with a set of unknown agents, with whom it can potentially establish new collaborations. Pseudocode for our interactive learning algorithm is shown in Algorithm 2. The algorithm takes as input the same arguments as the batch learning algorithm, except it includes an index specifying which agent in the stochastic game the learner is, and it does not

---

<sup>6</sup>A parameter controlling the linear combination of the team and bias function may be optimized as well.



---

**Algorithm 2** Interactive.Learning( $(I, S, A^I, T, R^I), k, \gamma, \mathcal{T}, B_\Theta$ )

---

**Require:** stochastic game  $(I, S, A^I, T, R^I)$ ; agent index  $k$ ; discount factor  $\gamma$ ; team function  $\mathcal{T}$ ; and parameterized bias function family  $B_\Theta$ .  
initialize policy  $\pi^k$  arbitrarily  
 $A^M := \times_i A^I$   $\triangleright$  Joint action set  
 $D := \{\}$   
**for** each round of play **do**  
    play round following  $\pi^k$   
    append new trajectory of play to  $D$   
     $R_\theta^S := \text{Batch\_Learning}((I, S, A^I, T, R^I), \gamma, D, \mathcal{T}, B_\Theta)$   
     $\forall s \in S, \pi^k(s) := a^k \in A^k$  where  $a \in \arg \max_a Q_{R_\theta^S}(s, a)$   
**end for**

---

include the batch of demonstrations. The agent begins by initializing its policy arbitrarily.<sup>7</sup> The first round of interaction produces a trajectory of joint behavior that the learner can now learn from. It does so by running our batch learning algorithm on the single trajectory of experience acquired thus far. It then proceeds to follow the learned policy in the next round. After the second round completes, the agent now has two trajectories on which it can run batch learning to update its social reward function and its behavior. This process then repeats for all rounds of interaction. Note that updating the policy after each round requires computing the  $Q$ -values under the newly learned social reward function. However, in practice, this computation can be performed lazily for each state the agent observes in the next interaction.

**Preliminary Results** Preliminary experiments suggest that both our batch and interactive collaborative learning algorithms can produce reasonable results in grid games. For batch learning, we hard code examples of agents following different equilibrium strategies, and show that our algorithm is able to recover the appropriate equilibrium when it is trained on demonstrations of behavior following that equilibrium. We also testing the agents on similar grid games after learning, to show that the learned behavior can generalize to other similar situations. For interactive learning, we show that agents can very quickly find an equilibrium, and that restarting learning from scratch can result in them finding a different one.

For both sets of experiments, we examine behavior in Hallway. The social reward function consists of a total welfare team reward function and a linear family of bias functions. The bias functions operate on 50 state-joint action binary features. The first 25 features are an indicator variable for selecting one of the 25 joint actions when the agents are in *conflict* and the latter 25 features are indicator functions for selecting one of the 25 joint actions when the agents are not a conflict. We define the agents to be in conflict when they are in the same row of the grid and each is closer to their opponent’s goal than they are to their own. Additionally, when one of the agents is one step away from their goal, none of the features activate, thereby preventing premature termination on the joint task. Since these features directly encode social rewards for joint actions, rather than higher-level goals, when training we use RH-MLIRL with a horizon of 1. We do not expect this set of features to be complete or capable of capturing all the possible behaviors that might emerge in a grid game. So an important area for future work is to identify features that are both expressive and generalize well. Nevertheless, the features used here provide a proof of concept for our approach.

We tested our batch learning algorithm on datasets in which agents followed three different equilibrium strategies. The first (Figure 2a) shows the orange player going north and then traveling along the top row of the grid, while the blue player goes south and then travels along the bottom row. The second (Figure 2b) shows the blue player going north and then traveling along the top row while the orange player waits, then goes east, and then waits again for blue to catch up so they can enter their goals together. The third (Figure 2c) has the orange player go south and then traveling along the bottom row, while the blue player immediately goes west toward its goal and then waits two steps.

For each of these three equilibria we generated five demonstrations, two of which described the behavior exactly, and the other three of which contained slight deviations (for example, moving back to the center row before reaching the end). We then separately trained our batch algorithm on each set of demonstrations.

---

<sup>7</sup>In practice, the policy can be initialized to the agent’s role in the joint policy that follows from the team function alone.

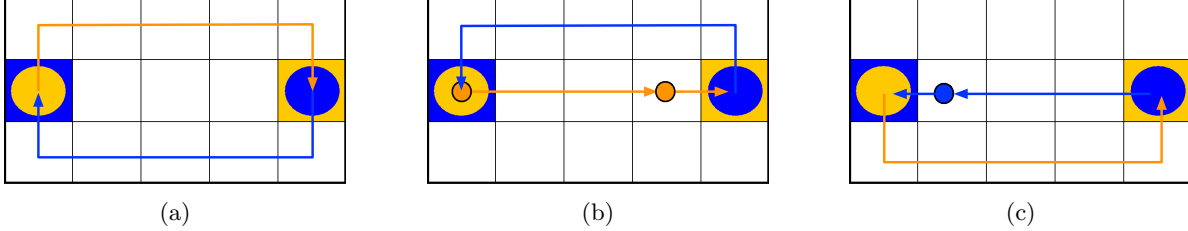


Figure 3: The three different equilibria the batch learning algorithm trained on in Hallway. Colored arrows indicate the path of the player of the same color. A small circle indicates a wait action.

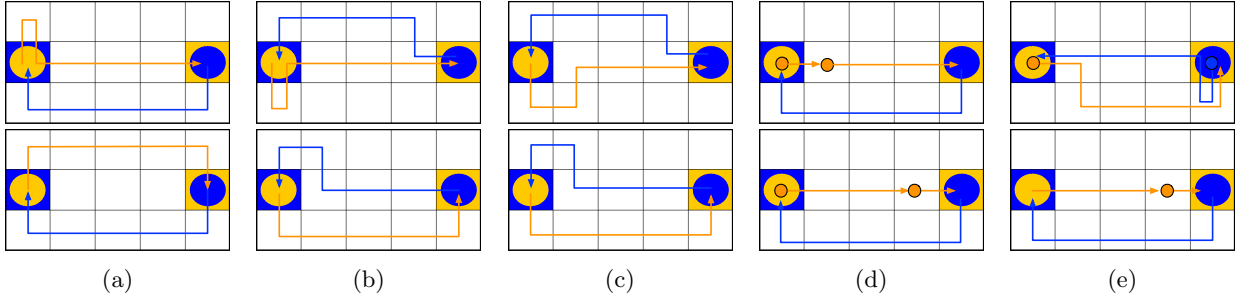


Figure 4: The learning results from each of the 5 interactive matches (a-e). The top image for a match is the first round of interaction of the agents. The bottom image for a match is the learned behavior. For matches b-e, the last move of an agent that reaches their goal sooner may be to wait for two steps, or to move north and/or south and back again; only one example is illustrated.

In all cases, the learned social reward function motivated behavior that consistently replicated the data on which it was trained. Additionally, we planned using the learned social reward function in a larger 7x3 hallway grid game, and in all cases, obtained the corresponding equilibrium behavior.

In our interactive learning experiment, we played two interactive learning agents against one another. The agents played five matches, each one consisting of five rounds. Each match began fresh, without any knowledge of previous matches, to see whether different equilibria could emerge. The social reward function was initialized to the total welfare joint policy. In each match, the agents very quickly (typically after the first round) converged to an equilibrium that persisted for the rest of the match. However, in each match, a different equilibrium was learned, showing that agents employing our algorithm quickly adapt to the particular shared experience they have with other agents.

Figure 3 shows the results of learning in each of the five matches (with the top image showing the behavior in the first round and the bottom image showing the learned behavior). While different equilibria emerged in different matches, once learned, that equilibrium was consistently employed in all subsequent rounds. Of particular interest are the results in the fifth match (Figure 3e). In the first round of this match, both players wait (a suboptimal joint action), and Blue then moves in an odd looping motion. Nevertheless, the agents learn from the experience, ultimately finding collaborative behavior that works well for them both.

### 3.3 Generalized Inverse Reinforcement Learning

To complement the new computational model we put forth, the new technology that we plan to develop as part of this proposal is a new model of IRL, and accompanying algorithms, which allows for both positive and negative (not necessarily binary) labels, and, moreover, which also allows for partial feedback, meaning subtrajectories (i.e., options [49]) can be ascribed such labels.

To this end, in this section, we propose a novel variant of the inverse reinforcement learning problem, which we call **Generalized IRL** (GIRL). Whereas the input to standard IRL is a data set of trajectories and the output, a reward function, the input to a GIRL problem, is a data set of *labelled* trajectories. In its

most basic formulation,<sup>8</sup> each label is either positive (+1) or negative (−1), indicating whether the associated trajectory is a good example of some desired behavior, or a bad example of some desired behavior. The goal is to recover a reward function that motivates the desired behavior.

The plate diagram that describes our GIRL problem formulation is shown in Figure 4. In this model, we interpret the label ( $L$ ) associated with a trajectory of size  $N$  to be some random function of what the evaluator thought about each of the action selections ( $A$ ) exhibited at each of the states ( $S$ ) in the trajectory. However, these step-wise evaluations ( $X$ ) are unobserved in the data.

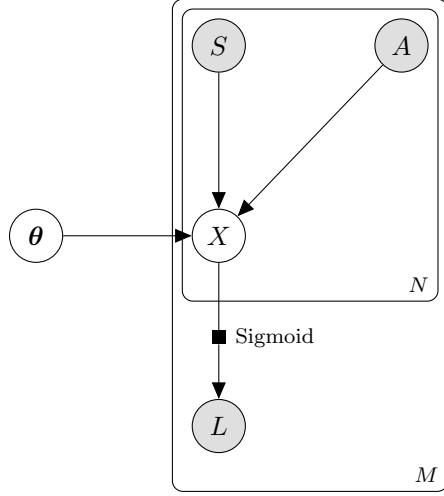


Figure 5: Generalized IRL

For illustrative purposes, we adopt the probabilistic maximum likelihood problem formulation, in which we seek a reward function that maximizes the likelihood of the data, but Bayesian reasoning could likewise be applied here. To find the reward function that maximizes the likelihood of the data, we first define a probability model for the system that is represented in Figure 4. We model the probability that an action is evaluated as good or not as proportional to its selection probability according a softmax policy computed for the reward function with parameters  $\theta$ . Specifically:

$$\Pr(x_i = +1 \mid s, a, \theta) = \pi(s, a \mid \theta) \quad (2)$$

$$\Pr(x_i = -1 \mid s, a, \theta) = 1 - \pi(s, a \mid \theta), \quad (3)$$

Here  $\pi(s, a \mid \theta)$  is the softmax policy over  $Q$ -values, assuming a reward function parameterized by  $\theta$ .

For the probability distribution of  $L$ , given the sequence of  $N$  step-wise labels, we would like a distribution that has the property that as more step-wise labels are positive, the probability of a positive trajectory label increases (and vice versa). Although there are many possible distributions that satisfy this property, for concreteness, we choose the sigmoid function. That is,

$$\Pr(L = +1 \mid X_1, \dots, X_n) = \frac{1}{1 + e^{-\phi \sum_i X_i}} \quad (4)$$

$$\Pr(L = -1 \mid X_1, \dots, X_n) = 1 - \Pr(L = +1 \mid X_1, \dots, X_n), \quad (5)$$

Here  $\phi$  is a parameter that controls how readily step-wise labels influence the trajectory labels. For example, when  $\phi = 0$ , trajectory labels are assigned uniformly at random, without any regard for of step-wise labels. As  $\phi \rightarrow \infty$ , the sigmoid converges to a step function in which a trajectory containing even one more positive step-wise label than negative will deterministically lead to a positive trajectory label (and vice versa).

**EM Algorithm** The problem with estimating the parameters ( $\theta$ ) of our reward function is that we have a latent variable vector  $X$  (or rather, some of the elements of the  $X$  vector are latent, and some may be observed), which prevents us from easily computing the likelihood of the model and maximizing its parameters. The EM approach to solving this problem is to first choose values for  $\theta$ ; then choose a new  $\theta$  that maximizes the expected value of the log likelihood function, where the distribution in the expectation is the probability of the latent variables (the  $X$ s in our case), given the observed data and previous choice of  $\theta$ . The maximization can be performed using gradient ascent, and this process repeats until convergence.

To simplify the exposition of the EM algorithm, we introduce the notation  $\mathbf{x}_k$  to indicate the values of the subset of observed (i.e., known) elements in an  $\mathbf{x}$  vector, and  $\mathbf{x}_u$  to represent a possible assignment to the subset of unobserved values. Using this notation, the EM algorithm is summarized in Algorithm 3:

The log likelihood given our plate diagram (i.e., the joint likelihood of the data and our parameters, given label  $l$  and a vector  $\mathbf{x}$ ) is  $\mathcal{L}(\mathbf{s}, \mathbf{a}, \theta \mid l, \mathbf{x}) = \Pr(l \mid \mathbf{x}) \prod_i \Pr(x_i \mid s_i, a_i, \theta)$ . Likewise, the log likelihood is  $\log \mathcal{L}(\mathbf{s}, \mathbf{a}, \theta \mid l, \mathbf{x}) = \log \Pr(l \mid \mathbf{x}) + \sum_i \log \Pr(x_i \mid s_i, a_i, \theta)$ . Additionally, the gradient of the log likelihood is  $\nabla_{\theta} \log \mathcal{L}(\mathbf{s}, \mathbf{a}, \theta \mid l, \mathbf{x}) = \sum_i \frac{\nabla_{\theta} \Pr(x_i \mid s_i, a_i, \theta)}{\Pr(x_i \mid s_i, a_i, \theta)}$ . This gradient forms the basis for the maximization step in EM.

<sup>8</sup>Labels need not be binary; they can be continuous-valued.

---

**Algorithm 3** EM Algorithm for GIRL
 

---

**Require:** initial  $\theta^0$ , and data  $s, a, x_k, l$   
**for**  $t = 0$  to  $T$  **do**  
      $\theta_{t+1} \leftarrow \arg \max_{\theta'} E_{x_u \sim \Pr(x_u | l, x_k, s, a, \theta)} [\log \mathcal{L}(s, a, \theta' | l, x)]$   
**end for**

---

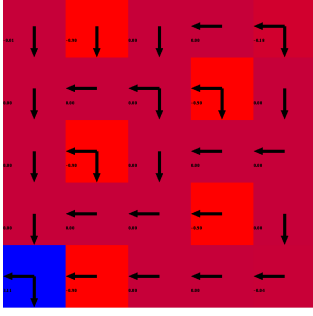


Figure 6: Two positive labelled trajectories.

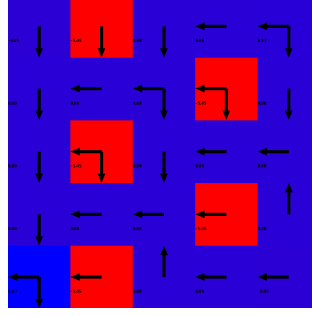


Figure 7: Positive & negative labelled trajectories.

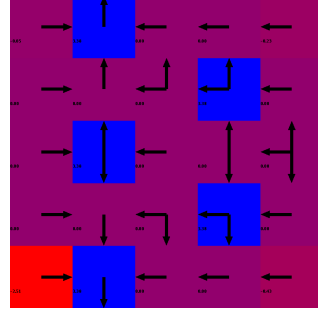


Figure 8: Negative & positive labelled trajectories.

The expectation in our EM algorithm is the expected value of the log likelihood under some candidate parameter  $\theta'$ , assuming unknown values are distributed according to  $\theta$ , is:

$$E_{x_u \sim \Pr(x_u | l, x_k, s, a, \theta)} [\log \mathcal{L}(s, a, \theta' | l, x)] = \sum_{x_u} \Pr(x_u | l, x_k, s, a, \theta) \log \mathcal{L}(s, a, \theta' | l, x_k, x_u).$$

Unfortunately, the number of assignments enumerated in the expectation's outer sum grows exponentially. A resolution to this problem is to estimate the expectation with sampling. However it is not easy to sample from  $\Pr(x_u | l, x_k, s, a, \theta)$ . On the other hand, it is easy to sample from  $\Pr(x_u | s, a, \theta)$  (removing the conditioning on the label). So, using importance sampling, we replace the expectation computation with the following sample-estimate, where  $x_u^j$  is the  $j$ th sample from the distribution  $\Pr(x_u | s, a, \theta)$ :

$$\frac{1}{C} \sum_{j=1}^C \frac{\Pr(x_u^j | l, x_k, s, a, \theta)}{\Pr(x_u^j | s, a, \theta)} \log \mathcal{L}(l, x_k, x_u^j | s, a, \theta) = \frac{1}{C} \sum_{j=1}^C \frac{\Pr(l | x_k, x_u)}{\Pr(l | x_k, s, a, \theta)} \log \mathcal{L}(l, x_k, x_u^j | s, a, \theta). \quad (6)$$

Lastly, while a straightforward computation of  $\Pr(l | x_k, s, a, \theta)$  requires marginalizing over all possible assignments to  $x_u$ , this term can be actually computed efficiently using dynamic programming, because the probability of a label is a sigmoid function that operates on the *sum* of the entries in the  $x$  vector. Putting all of these pieces together yields a tractable version of the EM algorithm for generalized IRL.

**Preliminary Experiments** We ran some (very) preliminary experiments to demonstrate the potential of an EM algorithm we developed for the GIRL framework.<sup>9</sup>

The GIRL agent is operating in a grid world, with puddles of multiple colors: red, blue, green, yellow and pink. The agent can take north, south, east and west actions. If there is a wall in the way, the agent does not move at all. In these experiments, we gave the GIRL agent two very short trajectories, and a corresponding pair of positive and negative labels. We will show that varying those labels leads the agent to infer different reward functions, and consequently, different behaviors.

In both trajectories, the agent begins in the middle cell of the bottom row of the grid. The first trajectory has the agent move two steps to the west, through the blue puddle, to reach the red puddle. The second trajectory has the agent move a single step west, into the blue puddle.

In the first experiment, we label both trajectories +1; in the second experiment, we label the first trajectory +1, and the second -1; and in the third experiment, we reverse these labels. The learned policies (along with corresponding, but illegible, rewards) are shown in Figures 5, 6, and 7, respectively.

---

<sup>9</sup>Because the problem domain in these preliminary experiments was small (a  $5 \times 5$  grid), we did not use importance sampling.

In experiment 1, we signal to the agent that both observed behaviors are good. Consequently, the agent learns to value the red puddle more than any other cell in the grid, and is intent on reaching it via a shortest path from everywhere else. In particular, the agent doesn’t learn to linger on the blue puddle, because in one of the trajectories walking over the blue to reach the red was labelled good.

In experiment 2, we again teach the agent that reaching the red puddle is good, but this time we also signal that walking into the blue puddle is bad. The agents again learns to get to the red puddle as quickly as possible, but this time it avoids all blue puddles. possible.

In experiment 3, we teach the agent that the red puddle is to be avoided, and getting to the blue puddle is good. This time, the agent prefers the blue puddles to the red, and it especially values the ones where it can remain in the blue cells indefinitely.

Thus, we have demonstrated that agents can learn utility functions within the GIRL framework using both positive and negative labels. We expect this capability to prove essential for humans hoping to teach agents to collaborate. Beyond what is possible now, with GIRL algorithms, we, the humans, will be able to critique agents during our interactions with them, by ascribing to their behaviors (entire trajectories and subtrajectories) both positive and negative feedback, based on which the agents will infer social reward functions, which we expect to generalize well within the learning domain and across learning environments.

## 4 Plan, Deliverables, and Evaluation

**Plan** The three-year timeline for our proposed project is as follows:

- **Year 1:** Demonstrate that a socially-rational agent can successfully learn collaborative behavior from batches of demonstrations (offline), and can generalize across related games. Carry out simulation experiments on machine-machine pairings, varying both the algorithms and the environments. Complete the development of an experimental test bed that pairs humans with artificial agents.
- **Year 2:** Establish that IRL algorithms within the GIRL framework can learn more accurately and more efficiently (i.e., from fewer data) than classic IRL algorithms, because of the ability of the teacher to offer both positive and negative feedback. Incorporate GIRL learning algorithms into the Year 1 experimental test bed, and show how collaborative behavior can emerge in real-time during human-machine interactions. Develop the proposed “Social autonomous driving” course (see Broader Impacts).
- **Year 3:** Extend as necessary, and then evaluate our approach in real-world applications, specifically to the “go fetch” scenario (see Broader Impacts).

**Deliverables** The proposed project includes the following deliverables:

1. An *experimental platform* in which humans can interact with artificial agents, thereby providing a space for artificial agents in general, and our socially-rational agents, specifically to demonstrate their capacity to represent, learn, and apply cooperative behaviors in context-appropriate ways.
2. *Software infrastructure for machine learning algorithm development and experimentation* that enable the representation and learning of cooperative behaviors in simulated stochastic games. Code will be made available as a part of the Brown-UMBC Reinforcement Learning and Planning library.<sup>10</sup>
3. Data comprised of the behaviors produced by both humans and machines on a test bed of stochastic games, varying the degrees of complexity and uncertainty, as well evaluations of those behaviors. It is our intent that these data will form a set of *computational benchmarks*, which other researchers will be free to use to evaluate their own approaches to collaborative learning.

---

<sup>10</sup>The open source Brown-UMBC Reinforcement Learning and Planning (BURLAP) library (<http://burlap.cs.brown.edu>) is the software infrastructure that will be used for algorithm development in this project. Developed by James MacGlashan under the direction of co-PI Littman, the system provides a flexible and powerful implementation of a wide variety of reinforcement-learning algorithms and environments, and has been extensively used in both published research and education (including in a reinforcement-learning MOOC). The algorithm which formed the basis of the work we describe in this proposal (RHC-MLIRL) is available as part of the BURLAP distribution.

4. Syllabus for an undergraduate class on “Social autonomous driving” in which students learn to endow autonomous vehicles with the ability to represent, learn, and apply cooperative driving algorithms.

In year two, we plan to organize a workshop (e.g., a AAAI symposium) on the machine learning of socially rational behavior. Such a gathering would be timely, as a number of subareas are grappling with related problems. We will include researchers in human–robot collaboration [18], intelligent software **\*\*Amy: missing reference\*\***, end-user programming of household devices [51], agent-based bidding [54], multi-agent reinforcement learning [48], computational social cognition [3], and others. **\*\*Amy: do you want to reference self-driving cars here?\*\***

**Evaluation** Before scaling to embodied robots, **\*\*Amy: Michael, are we going to do this? if not, tone down to whatever our real-world application plans are\*\*** we will evaluate how our socially-rational agents fare in collaborative learning tasks using an experimental platform that simulates interactions between humans and artificial agents. For example, we will compare how often and how quickly humans are able to arrive at coordinated plans among themselves, versus with a socially-rational agent. To test the expressiveness of the space of learnable solutions afforded by our approach, we will also measure tendencies among humans to converge to plans with certain characteristics, and whether those characteristics vary when there are socially-rational agents in the mix. We will also assess the quality of the plans under these varying conditions. Finally, subjective evaluations from the human participants will be used to assess whether interacting with agents feels natural.

For the batch learning algorithm, we will develop new quantitative metrics to assess behavioral similarity so that trajectories produced by our batch learners can be quantitatively compared to the trace data produced by human participants. Additionally, we plan to design a sort of “Turing test” for our problem domain, whereby we will have two humans interact with each other for a number of rounds and then, without notification, we will fork those interactions into two, where each human is now playing with a socially-rational agent that trained on the history of interactions between the two humans. At the end of all interactions, we will ask the user whether they noticed anything strange partway through the interactions, such as a change of partner. Their subjective response to whether they thought their partner switched at any point will be compared to their subjective responses in other treatments, such as: (1) no switching; (2) swapping their partner for another human; (3) swapping their partner for another human who watched the previous interactions; (4) swapping their partner for a socially-rational agent that did not perform any learning.

Finally, we are particularly interested in designing flexible learning algorithms that generalize at least as well as humans do to new state spaces. To this end, we will have participants play a series of related grid games, and then using evaluations similar to those listed above, we will test whether our learning agents can generalize to new games as well as other humans can. Likewise, **\*\*Amy: Michael, do we want to say something about generalizing to other players? maybe humans are better at this than our agents will be?\*\***

**\*\*Amy: Michael, how are we going to evaluate our approach in our proposed “real-world” (i.e., non-grid game) applications?\*\***

In year three we will test our approach in a real-world domain: robot-human collaboration. Even if (as we expect), our simulated worlds will suggest that our learning agents could collaborate better with people than traditional approaches, but the capabilities of robots may inhibit the effectiveness of our approach (e.g., robots might move slower than people expect). We will evaluate our algorithm in real-world scenarios by comparing its performance at controlling a robot to when another human controls the robot via teleoperation. We expect this to inspire novel problems tuned to real-world problems. For example, computation time may be a limiting factor to the success of our machine agents, and if so, we will develop more efficient implementations.

## 5 Broader Impacts

Recent trends in computer science and artificial intelligence are moving us toward increasing dependence on human–computer collaborative systems in which people and software make decisions that impact one another. Some fields, such as human–computer interaction (HCI), focus on systems in which a human being is primarily in control, and the computer’s decisions are assessed in terms of the positive impact they have

on the user. An interesting recent example is Centaur chess<sup>11</sup> in which a human and machine team up to play on the same side in chess, both making decisions, but with the machine acting as an advisor and the human deciding which moves to actually make. Other fields, like crowd sourcing human computation [52], combine human and machine expertise such that the machine is the ultimate arbiter of behavior and a group of human beings act as lower level computational components.

True collaboration, however, does not start with one participant being assigned to a leadership role. Instead, the various agents need to dynamically negotiate their roles and jockey for position, discovering when and how to trust each other to move forward. A concrete example is in the context of self-driving cars. Cars share the road with each other and must carefully choose when to be responsive to other vehicles and when to assert themselves to create a situation that benefits them. Doing so makes a significant difference in the driving’s effectiveness [?]. A related problem arises when a self-driving wheelchair attempts to move through a group of pedestrians [?]. More generally, robots that interact with people in the physical world need to navigate the complex give-and-take of establishing mutually beneficial behaviors where possible.

Indeed, a wide variety of human–machine interaction problems would be positively impacted by the technology we are proposing. As such, the project is synergistic with the main efforts of Brown University’s Humanity Centered Robotics Initiative (HCRI), of which co-PI Littman is co-director. Specifically, within HCRI, there are ongoing efforts to design robotic systems that interact with people and support independent living tasks (e.g., gerontechnological support for aging in place and robotic sous chefs). For an elderly person (or chef) to trust and collaborate on tasks with a machine effectively, the machine must act in a manner that the elderly person (or chef) expects. This project creates the foundation for these important applications.

The most direct application of our ideas would be in scenarios in which a robot is sharing a space with bystanders (human or robotic) who do not necessarily share its goals—a personal grocery-store assistant, a socially aware vacuum, an automated hospital linen cart, a package delivery drone. To further develop our ideas in a physical environment like these, and moreover in embodied robots, we plan to build a “go fetch” application on a mobile-manipulation platform that is scheduled for deployment in the Brown Computer Science department in **\*\*Amy: 2017??. can we stay some more convincing, so that the reviewers can be positive this environment will exist for us to work with\*\*** To be successful in this capacity, a robot will need to be able to represent, learn, and apply collaborative behavior in (at least) three separate spheres:

- navigating smoothly through the atrium and hallways, finding appropriate ways to avoid collisions, and adopting regularities that make it easy for people to coordinate their own movements with it;
- riding in the elevator, which brings with it some related but distinct expectations for coordinating (moving in and out of the doors, standing inside, asking for help with the buttons, deciding what to do when obstructed, etc.);
- understanding how to best signal a handoff of the object it is fetching, by interrupting people without being too timid or too boorish.

Interest in robotics and machine learning is growing within the Brown CS department. Presently, our faculty is designing a new robotics introductory course sequence that includes classes on basic robotics algorithms, hardware design, and human-centered evaluation of robotics systems. Inspired by the ideas in this proposal, we hope to contribute to this effort a new intermediate robotics course for undergraduates called “Social autonomous driving.” As a test platform, we will use MIT’s Duckietown platform<sup>12</sup>. **\*\*Amy: Michael, please insert footnote or comment about how this is already underway.\*\*** The emphasis of the class will be on developing robots capable of 1. maneuvering test vehicles in this test environment, and 2. interacting smoothly with other robots doing the same thing. The latter will be achieved by endowing the robots with socially-rational capabilities, thereby enabling them to adapt to local “customs” and driving styles (such as the “Pittsburgh left” or the “Boston rotary”).

The co-PIs have a strong record of mentoring students from underrepresented groups, and intend to continue pursuing these efforts in the context of this proposal. In particular, we plan to integrate our research

<sup>11</sup>[http://www.huffingtonpost.com/mike-cassidy/centaur-chess-shows-power\\_b\\_6383606.html](http://www.huffingtonpost.com/mike-cassidy/centaur-chess-shows-power_b_6383606.html)

<sup>12</sup><http://duckietown.mit.edu/>

on human-machine collaborations into Artemis, a free summer program directed by co-PI Greenwald that introduces rising 9th grade girls to computational thinking. For example, we might have the Artemis girls teach a robot to collaborate with them on routine tasks, such as navigation or object search.

Finally, we expect to publish the results of the proposed research in top-tier archival, conference proceedings and journals with high impact factors, and to present our work at innovative, non-archival workshops (e.g., the AAAI symposia).

## 6 Results from Prior NSF Support

**Intellectual Merit** Greenwald is currently funded by NSF (RI: Small-1217761, 2012-2016, \$450k) to build artificial agents that assist humans with decision making in information-rich and time-critical environments like online markets. This work is ongoing, but some preliminary publications include [17, 20, 21, 22]. Previously, she was funded to develop “Methods of Empirical Mechanism Design (EMD)” (CCF: Medium-0905234, 2009-2012, \$850k, with Mike Wellman). This project expanded the scope of mechanism design beyond the small-scale, stylized, or idealized domains most previous work was limited to. Publications include [34, 35, 55]. Before that, she received two prior NSF awards, a PECASE CAREER grant, “Computational Social Choice Theory” (IIS: Career-0133689, 2002-2007, \$375k), and “Efficient Link Analysis” (IIS: Small-0534586, 2005-2008, \$363k), which focused on ranking web pages and other social networks with hierarchical structure.

Littman is a co-PI on “RI: Medium: Collaborative Research: Teaching Computers to Follow Verbal Instructions” (No. 1065195, 9/11-8/14, \$704K) and “RI: Small: Collaborative Research: Speeding Up Learning through Modeling the Pragmatics of Training” (No. 1319305, 10/13-9/15, \$440k). These proposals developed autonomous agents that extract preferences from people via verbal interaction and reward feedback [30, 32, 33].

**Broader Impact** Broader impacts of our past work include undergraduate and graduate student training, design and evaluation of learning modules for outreach programs Learning Exchange and Artemis, and organizing two major computer science conferences (Littman) [14]. We also published novel benchmarks for evaluating grounded language learning programs (Littman), and developed a simulation platform for empirical game-theoretic analysis (Greenwald).

Greenwald also had an additional NSF award (EAGER: 1059570, 2010-2012, \$90k) which supported the expansion and evaluation of the Artemis Project, a free, summer program in which Brown computer science women teach computer science to rising 9th grade girls. This program has the dual effect of empowering Brown women, while at the same time, exposing girls to computer science. With this money, she successfully expanded Artemis to BU, and ran pilot programs at Columbia and UMBC.

Greenwald also has \$5k in funding from the Tides Foundation, through the IgniteCS program, for a project entitled Bringing Computer Science Education to Providence Public Schools.



## Data Management Plan

Our data management plan is designed to meet the letter and spirit of NSF data-management policy as described in the following documents:

- NSF 11-1 Grant Proposal Guide: Sec. II.C.2.J Special Information and Supplementary Documentation.
- NSF AAG11-001 Award and Admin. Guide Section VI.D.4.: Dissemination and Sharing of Research Results.

## Expected Data and Other Products

The expected data generated by this research includes open-source software libraries, robotic interaction data sets, and performance benchmarks. The software will be collected and stored in a GIT repository. The PIs will provide open-source implementations of parts of the software dealing with decision making and learning. A common GIT file repository will be employed for code development and data set storage. For each human experiment, the collected data will include:

- definition files with a description of the study materials used.
- trajectories with time-stamped states, controls, costs.

These materials will be completely anonymized with no personally identifying attributes. For machine experiments, we will build on the BURLAP library and will make our extensions available via GIT.

## Period of Data Retention

Data will be retained for a minimum of five years after completion of the grant.

## Data Formats

Human trajectory data will be stored in ASCII files along with a javascript program for interpreting/visualizing them. Other data will be stored in either ASCII or binary form.

## Data Access and Sharing

The software and collected log files will be made publicly available at the time of the first release of an associated publication. The data will be posted on a Brown-maintained web server, then moved to a Brown-maintained archive at the end of the project. All data is network accessible and routinely backed up onto physical drives. Upon substantial increase in data size, the PIs will consider data-management services such as available from Brown or Google Cloud. The PIs will specify conditions for data use, designate the original source of data, and associate it with the NSF grant number.

## References

- [1] Brenna D. Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, 57(5):469–483, 2009.
- [2] Monica Babes, Vukosi N. Marivate, Michael L. Littman, and Kaushik Subramanian. Apprenticeship learning about multiple intentions. In *International Conference on Machine Learning*, pages 897–904, 2011.
- [3] C. L. Baker and J. B. Tenenbaum. Modeling human plan recognition using Bayesian theory of mind. In G. Sukthankar, R. P. Goldman, C. Geib, D. Pynadath, and H. Bui, editors, *Plan, Activity and Intent Recognition: Theory and Practice*, pages 177–204. Morgan Kaufmann, 2014.
- [4] Andrew G. Barto, S. J. Bradtke, and Satinder P. Singh. Learning to act using real-time dynamic programming. *Artificial Intelligence*, 72(1):81–138, 1995.
- [5] Richard Bellman. *Dynamic Programming*. Princeton University Press, Princeton, NJ, 1957.
- [6] Dimitri P. Bertsekas. *Dynamic Programming: Deterministic and Stochastic Models*. Prentice-Hall, Englewood Cliffs, NJ, 1987.
- [7] M. Blanco, D. Englemann, and H. T. Normann. A within-subject analysis of other-regarding preferences. *Games and Economic Behavior*, 72:321–338, 2011.
- [8] Craig Boutilier, Thomas Dean, and Steve Hanks. Decision-theoretic planning: Structural assumptions and computational leverage. *Journal of Artificial Intelligence Research*, 11:1–94, 1999.
- [9] C. F. Camerer. *Behavioral Game Theory: Experiments in Strategic Interaction*. Princeton University Press, Princeton, NJ, 2003.
- [10] Jaedeug Choi and Kee-Eung Kim. Nonparametric bayesian inverse reinforcement learning for multiple reward functions. In *Proceedings of the 25th International Conference on Neural Information Processing Systems*, NIPS’12, pages 305–313, USA, 2012. Curran Associates Inc.
- [11] Gregg Collins and Louise Pryor. Planning under uncertainty: Some key issues. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1567–1573, 1995.
- [12] A. Cournot. *Recherches sur les Principes Mathématiques de la Théorie la Richesse*. Hachette, 1838.
- [13] Constantinos Daskalakis, Paul W Goldberg, and Christos H Papadimitriou. The complexity of computing a nash equilibrium. *SIAM Journal on Computing*, 39(1):195–259, 2009.
- [14] Marie desJardins and Michael L. Littman. AAAI-13 preface. In *AAAI 2013*, 2013.
- [15] A.M. Fink. Equilibrium in a stochastic  $n$ -person game. *Journal of Science in Hiroshima University*, 28:89–93, 1964.
- [16] R. Fisman, S. Kariv, and D. Markovits. Individual preferences for giving. *The American Economic Review*, 97(5):1858–1876, 2007.
- [17] T. Goff, A. Greenwald, E. Hilliard, W. Ketter, A. Loomis, and E. Sodomka. JACK: A Java auction configuration kit. In *Agent-Mediated Electronic Commerce. Designing Trading Strategies and Mechanisms for Electronic Markets.*, volume 136 of *Lecture Notes in Business Information Processing*, pages 45–60. Springer, 2013.
- [18] Nakul Gopalan and Stefanie Tellex. Modeling and solving human-robot collaborative tasks using POMDPs. In *Robotics: science and systems 2015: workshop on model learning for human-robot communication*, 2015.

- [19] A. Greenwald and K. Hall. Correlated  $Q$ -learning. In *Proceedings of the Twentieth International Conference on Machine Learning*, pages 242–249, 2003.
- [20] A. Greenwald, J. Li, and E. Sodomka. Approximating equilibria in sequential auctions with incomplete information and multi-unit demand. In *NIPS '12: Advances in Neural Information Processing Systems*, volume 25, pages 2330–2338, December 2012.
- [21] A. Greenwald, J. Li, E. Sodomka, and M. L. Littman. Solving for best-responses in extensive-form games using reinforcement learning methods. In *RLDM '13: Proceedings of 1st Multidisciplinary Conference Reinforcement Learning and Decision Making*, pages 116–120, October 2013.
- [22] A. Greenwald, E. Sodomka, E. Stix, J. Stix, and D. Storch. An empirical analysis of auctioneer profitability in Quibids penny auctions. In *Agent-Mediated Electronic Commerce. Designing Trading Strategies and Mechanisms for Electronic Markets.*, volume 197 of *Lecture Notes in Business Information Processing*, pages 56–68. Springer, 2014.
- [23] J. Hu and M. Wellman. Nash  $Q$ -learning for general-sum stochastic games. *Machine Learning Research*, 4:1039–1069, 2003.
- [24] D. Kahneman, P. Slovic, and A. Tversky, editors. *Judgment under uncertainty: Heuristics and biases*. Cambridge University Press, 1982.
- [25] Michael Kearns, Yishay Mansour, and Andrew Y. Ng. A sparse sampling algorithm for near-optimal planning in large Markov decision processes. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence (IJCAI-99)*, pages 1324–1331, 1999.
- [26] Levente Kocsis and Csaba Szepesvári. Bandit based Monte-carlo planning. In *Machine Learning: ECML 2006*, pages 282–293, 2006.
- [27] M. Littman. Friend or foe  $Q$ -learning in general-sum Markov games. In *Proceedings of Eighteenth International Conference on Machine Learning*, pages 322–328, June 2001.
- [28] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Proceedings of the eleventh international conference on machine learning*, volume 157, pages 157–163, 1994.
- [29] Michael L. Littman. Friend-or-foe  $Q$ -learning in general-sum games. In *Proceedings of the Eighteenth International Conference on Machine Learning*, pages 322–328. Morgan Kaufmann, 2001.
- [30] Robert Loftin, James MacGlashan, Michael L. Littman, Matthew E. Taylor, David L. Roberts, and Jeff Huang. A strategy-aware technique for learning behaviors from discrete human feedback. In *Proceedings of the Twenty-Eighth Association for the Advancement of Artificial Intelligence Conference (AAAI-14)*, 2014.
- [31] Manuel Lopes, Francisco Melo, and Luis Montesano. Active learning for reward estimation in inverse reinforcement learning. In *Machine Learning and Knowledge Discovery in Databases*, pages 31–46. Springer, 2009.
- [32] J. MacGlashan, M. Babes-Vroman, M. desJardins, M. Littman, S. Muresan, S. Squire, S. Tellex, D. Arumugam, and L. Yang. Grounding english commands to reward functions. In *Proceedings of Robotics: Science and Systems*, 2015.
- [33] James MacGlashan and Michael L. Littman. Between imitation and intention learning. In *Proceedings of IJCAI*, 2015.
- [34] B. Mayer, E. Sodomka, and A. Greenwald. The price of independence in simultaneous auctions. In *AAMAS '13: Proceedings of the 12th International Conference on Autonomous Agents and Multiagent Systems*, pages 1227–1228, May 2013.

- [35] B. Mayer, E. Sodomka, A. Greenwald, and M.P. Wellman. Accounting for price dependencies in simultaneous sealed-bid auctions. In *EC '13: Proceedings of the 14th ACM Conference on Electronic Commerce*, pages 679–696, June 2013.
- [36] Bernard Michini and Jonathan P. How. *Bayesian Nonparametric Inverse Reinforcement Learning*, pages 148–163. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.
- [37] Andrew Y. Ng, Daishi Harada, and Stuart J. Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Proceedings of the Sixteenth International Conference on Machine Learning, ICML '99*, pages 278–287, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.
- [38] Andrew Y. Ng and Stuart Russell. Algorithms for inverse reinforcement learning. In *International Conference on Machine Learning*, pages 663–670, 2000.
- [39] Dean A. Pomerleau. *Neural network perception for mobile robot guidance*. Kluwer Academic Publishing, 1993.
- [40] Qifeng Qiao and Peter A. Beling. Inverse reinforcement learning with gaussian process. *CoRR*, abs/1208.2112, 2012.
- [41] Deepak Ramachandran and Eyal Amir. Bayesian Inverse Reinforcement Learning. *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, 51:2586–2591, 2007.
- [42] Tummalapalli Sudhamsh Reddy, Vamsikrishna Gopikrishna, Gergely Zaruba, and Manfred Huber. Inverse reinforcement learning for decentralized non-cooperative multiagent systems. In *Systems, Man, and Cybernetics (SMC), 2012 IEEE International Conference on*, pages 1930–1935. IEEE, 2012.
- [43] Peter C. Reiss and Frank A. Wolak. Structural econometric modeling: Rationales and examples from industrial organization. In J.J. Heckman and E.E. Leamer, editors, *Handbook of Econometrics*, volume 6A, chapter 64. Elsevier, 1 edition, 2007.
- [44] T.C. Schelling. *The Strategy of Conflict*. Harvard University Press, 1980.
- [45] John R. Searle. *The construction of social reality*. Simon and Schuster, 1995.
- [46] L.S. Shapley. Stochastic games. *Proceedings of the National Academy of Sciences of the United States of America*, 39:1095–1100, 1953.
- [47] Satinder Singh, Richard L Lewis, and Andrew G Barto. Where do rewards come from. In *Proceedings of the annual conference of the cognitive science society*, pages 2601–2606, 2009.
- [48] Eric Sodomka, Elizabeth Hilliard, Michael Littman, and Amy Greenwald. Coco-Q: Learning in stochastic games with side payments. *JMLR Workshop and Conference Proceedings: Proceedings of The 30th International Conference on Machine Learning*, 28(3):1471–1479, 2013.
- [49] Richard S. Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1–2):181–211, 1999.
- [50] Michael Tomasello, Malinda Carpenter, Josep Call, Tanya Behne, and Henrike Moll. Understanding and sharing intentions: The origins of cultural cognition. *Behavioral and brain sciences*, 28(05):675–691, 2005.
- [51] Blase Ur, Elyse McManus, Melwyn Pak Yong Ho, and Michael L. Littman. Practical trigger-action programming in the smart home. In *Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems (CHI)*, 2014.
- [52] Luis Von Ahn. Human computation. In *Design Automation Conference, 2009. DAC'09. 46th ACM/IEEE*, pages 418–419. IEEE, 2009.

- [53] Payne John W., Bettman James R., and Johnson Eric J. *The adaptive decision maker*. Cambridge University Press, 1993.
- [54] M. Wellman, A. Greenwald, and P. Stone. *Autonomous Bidding Agents: Strategies and Lessons from the Trading Agent Competition*. Intelligent Robotics and Autonomous Agents. MIT Press, Cambridge, 2007.
- [55] M. P. Wellman, E. Sodomka, and A. Greenwald. Self-confirming price prediction strategies for simultaneous one-shot auctions. In *UAI '12: Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence*, pages 893–902, August 2012.
- [56] Brian Ziebart, Andrew L. Maas, J. Andrew Bagnell, and Anind K. Dey. Maximum entropy inverse reinforcement learning. In Dieter Fox and Carla P. Gomes, editors, *AAAI*, pages 1433–1438. AAAI Press, 2008.
- [57] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. In *AAAI*, pages 1433–1438, 2008.
- [58] M. Zinkevich, A. Greenwald, and M. Littman. Cyclic equilibria in Markov games. In *Advances in Neural Information Processing Systems 18*, page To Appear. MIT Press, 2006.