

Learning Human Rewards by Inferring Their Latent Intelligence Levels in Multi-Agent Games: A Theory-of-Mind Approach with Application to Driving Data

Ran Tian, Masayoshi Tomizuka, and Liting Sun

Abstract—Reward function, as an incentive representation that recognizes humans’ agency and rationalizes humans’ actions, is particularly appealing for modeling human behavior in human-robot interaction. Inverse Reinforcement Learning is an effective way to retrieve reward functions from demonstrations. However, it has always been challenging when applying it to multi-agent settings since the mutual influence between agents has to be appropriately modeled. To tackle this challenge, previous work either exploits equilibrium solution concepts by assuming humans as perfectly rational optimizers with unbounded intelligence or pre-assigns humans’ interaction strategies *a priori*. In this work, we advocate that humans are bounded rational and have different intelligence levels when reasoning about others’ decision-making process, and such an inherent and latent characteristic should be accounted for in reward learning algorithms. Hence, we exploit such insights from Theory-of-Mind and propose a new multi-agent Inverse Reinforcement Learning framework that reasons about humans’ latent intelligence levels during learning. We validate our approach in both zero-sum and general-sum games with synthetic agents and illustrate a practical application to learning human drivers’ reward functions from real driving data. We compare our approach with two baseline algorithms. The results show that by reasoning about humans’ latent intelligence levels, the proposed approach has more flexibility and capability to retrieve reward functions that explain humans’ driving behaviors better.

I. INTRODUCTION

Our society is rapidly advancing towards robots that collaborate with and assist humans in daily tasks. To assure safety and efficiency, such robots need to understand and predict human behaviors. Modeling humans as reward-driven agents is particularly appealing since reward functions recognize humans’ agency and rationalize their actions. Inverse Reinforcement Learning (IRL) has been proved to be an effective approach to learning reward functions from human demonstrations. However, most previous IRL work is restricted to single-agent settings [1]–[5].

Learning reward functions in multi-agent settings is challenging, as it requires an interaction model that characterizes the mutual influence or the closed-loop dynamics among agents. Theory-of-Mind [6] is a mechanism that explicitly models humans’ beliefs over other humans’ decision-making process. Based on Theory-of-Mind, when a human interacts with his/her opponent, his/her cognitive reasoning process is nested and can be structured in a recursive fashion: “I believe

Ran Tian, Masayoshi Tomizuka, and Liting Sun are with University of California, Berkeley. ({rantian, tomizuka, litingsun}@berkeley.edu).

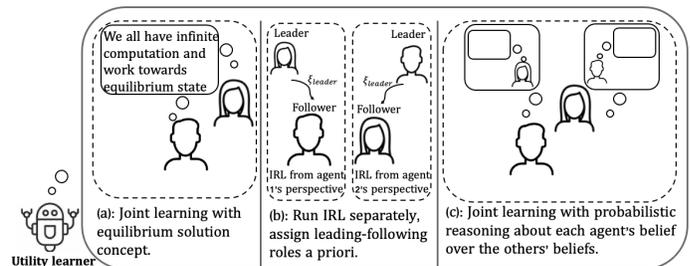


Fig. 1. Different IRL formulations in multi-agent interactions, (a): IRL approaches that exploits equilibrium solution concepts which assume agents have infinite intelligence levels. (b): IRL approaches that assign leader-follower roles *a priori*. (c): our approach that reasons about agents’ intelligence levels during reward learning.

that you believe that I believe...” We refer to the depth of such recursion as a human’s *intelligence level*, and it characterizes a human’s cognitive reasoning ability. Previous work on multi-agent IRL has leveraged equilibrium solution concepts of Markov Games to model interactions [7], [8] as shown in Fig. 1(a). Equilibrium solution-based approaches assume that humans have unlimited computation and infinite intelligence levels when making decisions. However, in real life, humans rarely pursue equilibrium solutions [9]. Hence, to account for such facts, non-equilibrium solution-based approaches have also been explored by many researchers in IRL algorithms. For instance, [10]–[13] exploited a leader-follower game to account for humans’ bounded intelligence in two-agent IRL settings. Instead of jointly learning agents’ reward functions, these approaches run IRL from each agent’s perspective with pre-assigned leader/follower roles (refer to Section II for more details).

In this work, we advocate that humans have latent intelligence levels that are naturally bounded and heterogeneous, and such an aspect should be accounted for in multi-agent IRL formulations. Specifically, we exploit insights from Theory-of-Mind to explicitly model humans’ interactions under bounded intelligence, and propose a multi-agent IRL formulation that reasons about humans’ heterogeneous intelligence levels during learning. Our key insight is: *reasoning about humans’ hidden intelligence levels (a type of latent cognitive states) adds more flexibility to multi-agent IRL algorithms and helps learn reward functions that explain and reconstruct human behaviors better.*

Overall, we make the following contributions towards multi-agent IRL:

Developing a framework for joint reward functions

learning in multi-agent settings without assumptions about perfect rationality or leader-follower roles. We exploit insights from Theory-of-Mind to account for humans’ bounded intelligence and extend the Maximum Entropy Inverse Reinforcement Learning (MaxEnt IRL) framework in multi-agent settings with explicit reasoning about humans’ latent intelligence levels during learning. We validate our approach in both zero-sum and general-sum games with synthetic agents, then apply our approach to learn human driver reward functions from real-world driving data.

Analyzing the advantages of reasoning about humans’ intelligence levels in multi-agent IRL with real driving data. We conduct, by our knowledge, the first quantitative comparison between two types of IRL algorithms using real driving data: the IRL algorithm that pre-assigns humans’ leader-follower roles during learning [10]–[13] and our approach that leverages the cognitive reasoning structure to reason about humans’ latent intelligence levels. We show that our approach provides more flexibility during learning and helps explain and reconstruct real human drivers’ behaviors better.

II. RELATED WORKS

Opponent modeling. Our work is an instance of opponent modeling. Traditional approaches in opponent modeling often model opponents’ behaviors based on past experience. In particular, [14]–[16] utilized type-based reasoning approaches that assume an opponent agent is one of known types and identify opponents’ types online based on past interactions. Other work also modeled opponents in terms of leader-follower relationship: ego agent models its opponents’ responsive actions as a probability distribution conditioned on its own action [17], [18]. Insights from Theory-of-Mind have been exploited to not only model opponents’ conditional distribution of actions but also model opponents’ beliefs over other agents’ beliefs, such nested beliefs (“I believe that you believe that I believe...”) can be explicitly modeled by approaches based on recursive reasoning [19]–[23]. Our approach to modeling opponents is based on the quantal level- k model that belongs to the recursive reasoning paradigm. In contrast to [22], we do not assume all agents perform depth-1 recursions but recognize the heterogeneity in agents’ recursive depths. In contrast to [23], we explicitly account for opponents’ sub-optimal behaviors.

Learning human driver reward functions. Many previous works in learning human drivers’ reward functions are extensions of MaxEnt IRL [4] and are restricted to single-agent settings [24]–[26]. Recently, the concept of quantal best response equilibrium (QRE) is exploited to extend MaxEnt IRL to multi-agent games. Equilibrium solution concepts assume all agents have infinite intelligence levels and have common knowledge of this. However, mounting evidence suggests that human behaviors often deviate from equilibrium behaviors in systematically biased ways due to their compromised / bounded intelligence [9], [27], [28]. Other approaches have extended MaxEnt IRL in two-agent settings by utilizing a

pre-assigned leader-follower relationship to model interactions [10]–[13]. Consequently, instead of jointly learning reward functions, they have to conduct MaxEnt IRL from each agent’s perspective separately. Moreover, in each separate IRL formulation, the interaction is simplified as an open-loop leader-follower game, where the opponent’s ground-truth trajectory is assumed to be accessible by the ego agent (Fig. 1 (b)). Therefore, no real interaction has been accounted for in the reward functions learning process. Furthermore, assigning the leader-follower relationship *a priori* requires careful selection for demonstrations because demonstrations that don’t align with the role assignment may lead to biased learning, i.e., incorrect human driver reward functions might be retrieved as demonstrated in VII.

Comparison between our approach and previous work. Our approach is different from the approaches in [7], [8] since we relaxed the assumptions about the equilibrium solution concept by leveraging recursive reasoning and recognizing humans’ bounded intelligence during learning. Our work is also distinguished from [29] by extending the learning algorithm to multi-agent settings and illustrating a practical application to learning human driver reward functions from real traffic data. In contrast to [10]–[13], our approach jointly learns human drivers’ reward functions in multi-agent games without assumptions about agents’ roles (Fig. 1 (c)).

III. PROBLEM FORMULATION

Multi-agent interaction formalization. We model the interactions between n agents ($n \in \mathbb{N}^+$) as a stochastic game defined by tuple $\langle \mathcal{S}, \mathcal{A}_1, \dots, \mathcal{A}_n, \mathcal{R}^1, \dots, \mathcal{R}^n, f, \gamma \rangle$, where \mathcal{S} denotes the finite state space, \mathcal{A}_i denotes the finite action space of agent i ($i \in \mathcal{P} = \{1, \dots, n\}$), $\mathcal{R}^i : \mathcal{S} \rightarrow \mathbb{R}$ denotes the reward function of agent i , $f : \mathcal{S} \times \mathcal{A}_1 \times \dots \times \mathcal{A}_n \rightarrow \mathcal{S}$ denotes the open-loop dynamics of the game, and γ is the discount factor. For agent i , we let $\pi_i : \mathcal{S} \rightarrow \mathcal{A}_i$ denote its deterministic policy. At each discrete time step t , agent i aims to maximize its expected total reward in state $s_t \in \mathcal{S} : \pi_i^* = \arg \max_{\pi} V_i^{\pi}(s_t)$, where $V_i^{\pi}(s_t) = \mathbb{E}_{a_{t:\infty}^{-i}} [\sum_{\tau=0}^{\infty} \gamma^{\tau} \mathcal{R}^i(s_{t+\tau}) | \pi, f]$ is the value function representing agent i ’s expected return starting from s_t , subject to its policy and the dynamics function. The expectation is taken with respect to the possible actions from other agents in the game, denoted by $a_t^{-i} = (a_t^j)$ with $j \in \mathcal{P} \setminus \{i\}$. Different reasoning strategies of the other agents will yield significantly different distribution for a_t^{-i} and impact the closed-loop dynamics of the game.

Reward learning as an optimization problem. We consider a reward-free stochastic game. We assume that agent i ’s reward function is parameterized by a collection parameter ω_i and $\mathcal{R}_{\omega_i}^i$ is differentiable with respect to ω_i . Such a reward function representation could be a neural network or a linear combination of a given set of features. Suppose we are given a demonstration set \mathcal{D} that contains multiple groups of interaction trajectories among n agents. Our objective is to jointly learn the reward functions for all agents that best rationalize the demonstrations. Namely, we want to find an optimal $\bar{\omega} = (\omega_1, \dots, \omega_n)$ that maximizes the likelihood of

observed demonstrations. We optimize for $\bar{\omega}$ in the following maximum likelihood estimation problem:

$$\max_{\bar{\omega}} \sum_{\xi \in \mathcal{D}} \log \mathbb{P}(\xi | \bar{\omega}) = \max_{\bar{\omega}} \sum_{\xi \in \mathcal{D}} \log \prod_{t=0}^{t_f} \mathbb{P}(\bar{a}_t | s_t, \bar{\omega}), \quad (1)$$

where $\bar{a}_t = (a_t^1, \dots, a_t^n)$ is the collection of all agents' actions at time step t , $\xi = \{(s_0, \bar{a}_0), \dots, (s_{t_f}, \bar{a}_{t_f})\}$ is an interaction trajectory, and $\mathbb{P}(\bar{a}_t | s_t, \bar{\omega})$ denotes the conditional probability of agents' actions given s_t and the current estimate of $\bar{\omega}$. Note that without specifications of the agents' abilities in reasoning about others' reasoning processes, it is hard to further decompose $\mathbb{P}(\bar{a}_t | s_t, \bar{\omega})$ to be related to each agent's policy π^i . Hence, in this work, we explicitly consider agents' bounded intelligence levels that characterize their reasoning abilities and exploit insights from Theory-of-Mind to formulate the multi-agent reward learning problem.

IV. HUMAN COGNITIVE REASONING MODELLING

Our instance of Theory-of-Mind follows from a recursive reasoning model - quantal level- k model (ql- k) [30] which is shown to be the state-of-the-art [31] in predicting human strategic behaviors. We briefly introduce the model below.

Quantal level- k model. Ql- k model represents humans' recursive reasoning structure in an iterative fashion, starting from ql-0 agents who are myopic agents with the lowest intelligence level. Then, a ql- k agent (agent with level- k intelligence), $k \in \mathbb{N}^+$, believes other agents as ql- $(k-1)$ agents (i.e., following ql- $(k-1)$ policies), and responds accordingly based on such beliefs. Note that with the ql-0 model as a base model, the ql- k policies are defined for every agent $i \in \mathcal{P}$ and for every intelligence level $k = 1, \dots, k_{\max}$ through a sequential and iterative process. We illustrate this process from agent i 's perspective in the left block of Fig. 2.

Specifically, given an initial state $s_t \in \mathcal{S}$, a ql- k agent i maximizes the following objective: $\max_{\pi^i} V^{i,k}(s_t)$, where $V^{i,k}(s_t) = \mathbb{E}_{\pi^{-i,k-1}} \left[\sum_{\tau=0}^{\infty} \gamma^\tau \mathcal{R}^i(s_{t+\tau}) | \pi^i, f \right]$ is the ql- k value function of agent i and $\pi^{-i,k-1}$ are the predicted ql- $(k-1)$ policies of agents $-i = \mathcal{P} \setminus \{i\}$ from agent i 's belief space. It follows that the optimal value function satisfies the following Bellman equation: $V^{*,i,k}(s) = \max_{a^i \in \mathcal{A}_i} \mathbb{E}_{\pi^{-i,k-1}} \left[\mathcal{R}^i(s') + \gamma V^{*,i,k}(s') | s' = f(s, a^i, a^{-i}), a^{-i} \sim \pi^{-i,k-1} \right]$. Then, we define the Q -value function as:

$$\begin{aligned} Q^{*,i,k}(s, a^i) &= \mathbb{E}_{\pi^{-i,k-1}} \left[\mathcal{R}^i(s) + \gamma V^{*,i,k}(s') \right] \\ &= \mathbb{E}_{\pi^{-i,k-1}} \left[\mathcal{R}^i(s) + \gamma \max_{a^{i'}} Q^{*,i,k}(s', a^{i'}) \right]. \end{aligned} \quad (2)$$

Note that $Q^{*,i,k}$ is in a form of Bellman equation and can be determined via value iteration. Then we define agent i 's ql- k policy using the following quantal best response function:

$$\pi^{i,k}(s, a^i) = \frac{\exp(\lambda^i Q^{*,i,k}(s, a^i))}{\sum_{a' \in \mathcal{A}_i} \exp(\lambda^i Q^{*,i,k}(s, a'))}, \quad (3)$$

where $\lambda^i \in (0, 1]$ is the rationality coefficient that controls the degree of agent i conforming to optimal behaviors. Start from the base case, ql-0 policies, by sequentially and iteratively solving for $\pi^{i,k}(s, a^i)$ via Eq. (2) and Eq. (3), we can compute ql- k policies for an arbitrary $k = 1, 2, \dots$ and for every agent. We note that when agent i predicts its opponents' ql- $(k-1)$ policies, it assumes that its opponents also use rationality coefficient λ^i to establish their ql- $(k-1)$ policies.

Summary. Thus far, we have introduced our instance of Theory-of-Mind model exploited to model humans' heterogeneous intelligence levels during strategic interactions and demonstrated how to represent humans' beliefs over other humans' beliefs in an iterative and sequential procedure.

V. COGNITION-AWARE MULTI-AGENT INVERSE REINFORCEMENT LEARNING

In this section, we incorporate the ql- k model into our multi-agent IRL formulation.

Learning with inference on latent intelligence levels. Though we have developed models that represent humans' cognitive reasoning under different intelligence levels, it is still difficult to further decouple $\mathbb{P}(\bar{a}_t | s_t, \bar{\omega})$ in (1) because humans' intelligence levels are latent cognitive states and can not be directly observed. We assume that humans' true intelligence levels are constant variables in a demonstration ξ . We let $\bar{k} = (k^1, \dots, k^n)$ denote the collection of humans' true intelligence levels in ξ and let \mathbb{K} denote a set that contains all possible intelligence levels ($k^i \in \mathbb{K}$). Then, as an approximation, we aim to maximize the expected log-likelihood of the demonstrations conditioned on estimates of the reward parameters $\bar{\omega}$ and a sample \bar{k} with respect to the current conditional distribution of \bar{k} given ξ and the current estimates of $\bar{\omega}$:

$$\begin{aligned} &\max_{\bar{\omega}} \sum_{\xi \in \mathcal{D}} \log \mathbb{P}(\xi | \bar{\omega}) \\ &\approx \max_{\bar{\omega}} \sum_{\xi \in \mathcal{D}} \sum_{\bar{k} \in \mathbb{K}^n} \log (\mathbb{P}(\xi | \bar{\omega}, \bar{k})) \mathbb{P}(\bar{k} | \xi, \bar{\omega}), \\ \mathbb{P}(\xi | \bar{\omega}, \bar{k}) &= \prod_{t=0}^{t_f} \prod_{i=1}^n \pi_{\bar{\omega}}^{i,k^i}(s_t, a_t^i), \mathbb{P}(\bar{k} | \xi, \bar{\omega}) = \prod_{i=1}^n \mathbb{P}(k^i | \xi, \bar{\omega}), \end{aligned} \quad (4)$$

where $\pi_{\bar{\omega}}^{i,k^i}$ denotes the ql- k^i policy of agent i and $\mathbb{P}(k^i | \xi, \bar{\omega})$ denotes the estimated probability distribution of agent i 's intelligence level conditioned on the current demonstration and estimate of reward parameters, which can be computed by applying recursive Bayesian inference along ξ :

$$\mathbb{P}(k^i = k | \xi_t, \bar{\omega}) = \frac{\pi_{\bar{\omega}}^{i,k}(s_t, a_t^i) \mathbb{P}(k^i = k | \xi_{t-1}, \bar{\omega})}{\sum_{k' \in \mathbb{K}} \pi_{\bar{\omega}}^{i,k'}(s_t, a_t^i) \mathbb{P}(k^i = k' | \xi_{t-1}, \bar{\omega})}, \quad (5)$$

where $\xi_t = \{(s_0, \bar{a}_0), \dots, (s_t, \bar{a}_t)\}$ ($t \leq t_f$) and the initial prior distribution is a uniform distribution over \mathbb{K} . Then, we set $\mathbb{P}(k^i = k | \xi, \bar{\omega}) = \mathbb{P}(k^i = k | \xi_{t_f}, \bar{\omega})$.

In summary, in order to evaluate the likelihood of a demonstration induced by the current estimate of agents' reward parameters using the ql- k model with the latent variable

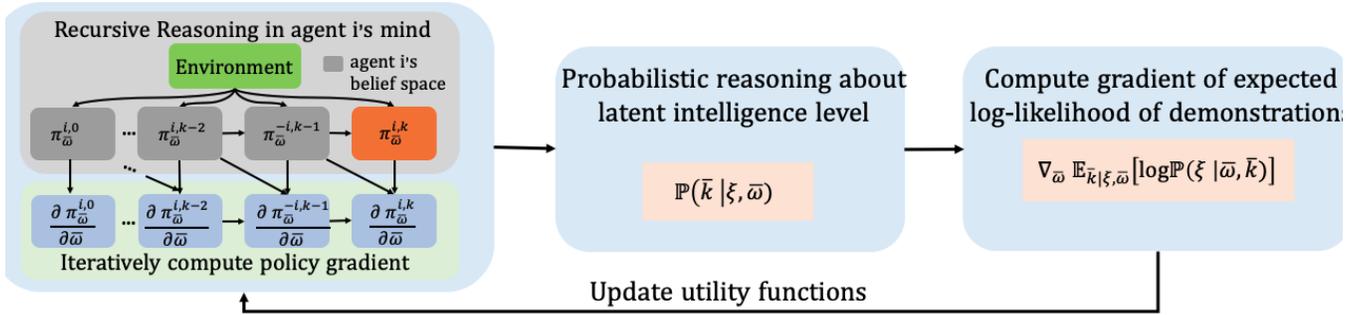


Fig. 2. Our proposed approach to learning reward functions in multi-agent settings. Left block: procedure for computing ql- k policies and gradients of ql- k policies with respect to weights of reward functions. Middle block: the computed $\pi_{\bar{\omega}}^{i,k}$ and $\frac{\partial \pi_{\bar{\omega}}^{i,k}}{\partial \bar{\omega}}$ for each i and k are used to compute the current conditional distribution of \bar{k} given ξ and the current estimate of $\bar{\omega}$. Right block: the information from the previous two blocks are exploited to compute the gradient of the expected likelihood of the demonstrations.

\bar{k} , we first compute the current conditional distribution of \bar{k} given ξ and the current estimate of agents' reward parameters $\bar{\omega}$ ($\mathbb{P}(k^i | \xi, \bar{\omega})$), then compute the expected likelihood of the demonstrations with respect to $\mathbb{P}(k^i | \xi, \bar{\omega})$. After that, the gradient of the objective function Eq. (4) can be used to find a locally optimal $\bar{\omega}$. Such a learning procedure is illustrated in Fig. 2.

Analytic Q-value gradient approximation. It follows from Eq. (4) that the gradient of the learning objective function with respect to $\bar{\omega}$ depends on the gradients of $\pi_{\bar{\omega}}^{i,k}$ and $\mathbb{P}(k | \xi, \bar{\omega})$ with respect to $\bar{\omega}$, which both further depend on the gradient of $Q_{\bar{\omega}}^{i,k}$ with respect to $\bar{\omega}$ according to Eq. (3). Due to the non-differentiable max operator associated with the computation of $Q_{\bar{\omega}}^{i,k}$ in Eq. (2), we utilize a smooth approximation of the Q value:

$$Q_{\bar{\omega}}^{i,k}(s, a^i) \approx \sum_{a^{-i}} \mathbb{P}(a^{-i} | s, \bar{\omega}) \left(\mathcal{R}_{\bar{\omega}}^i(s') + \gamma \left(\sum_{a^{i'}} (Q_{\bar{\omega}}^{i,k}(s', a^{i'}))^{\kappa} \right)^{\frac{1}{\kappa}} \right), \quad (6)$$

where the parameter κ controls the approximation error, and when $\kappa \rightarrow \infty$, the approximation becomes exact.

Let us assume that we have access to $\frac{\pi_{\bar{\omega}}^{i,k-1}}{\partial \bar{\omega}}$ of each agent, then the gradient of the value function $Q_{\bar{\omega}}^{i,k}$ with respect to $\bar{\omega}$ can be approximated as follows:

$$\begin{aligned} \frac{Q_{\bar{\omega}}^{i,k}}{\partial \bar{\omega}}(s, a^i) &\approx \sum_{a^{-i}} \left[\frac{\partial \Pi^{-i,k-1}}{\partial \bar{\omega}}(a^{-i} | s, \bar{\omega}) \left(\mathcal{R}_{\bar{\omega}}^i(s') \right. \right. \\ &\quad \left. \left. + \gamma \left(\sum_{a^{i'}} (Q_{\bar{\omega}}^{i,k}(s', a^{i'}))^{\kappa} \right)^{\frac{1}{\kappa}} \right) \right. \\ &\quad \left. + \mathbb{P}(a^{-i} | s, \bar{\omega}) \left(\frac{\partial \mathcal{R}_{\bar{\omega}}^i}{\partial \bar{\omega}}(s') + \gamma \frac{1}{\kappa} \left(\sum_{a^{i'}} (Q_{\bar{\omega}}^{i,k}(s', a^{i'}))^{\kappa} \right)^{\frac{1-\kappa}{\kappa}} \right. \right. \\ &\quad \left. \left. \cdot \sum_{a^i} \kappa \left(Q_{\bar{\omega}}^{i,k}(s', a^i) \right)^{\kappa-1} \frac{Q_{\bar{\omega}}^{i,k}}{\partial \bar{\omega}}(s', a^i) \right) \right], \quad (7) \end{aligned}$$

where $\frac{\partial \Pi^{-i,k-1}}{\partial \bar{\omega}}(a^{-i} | s, \bar{\omega})$ can be computed as follows:

$$\begin{aligned} &\frac{\partial \Pi^{-i,k-1}}{\partial \bar{\omega}}(a^{-i} | s, \bar{\omega}) \\ &= \sum_{j \in \mathcal{P} \setminus \{i\}} \frac{\pi_{\bar{\omega}}^{j,k-1}}{\partial \bar{\omega}}(s, a^j | \bar{\omega}) \prod_{e \in \mathcal{P} \setminus \{i,j\}} \pi_{\bar{\omega}}^{e,k-1}(s, a^e | \bar{\omega}), \quad (8) \end{aligned}$$

where $\frac{\pi_{\bar{\omega}}^{j,k-1}}{\partial \bar{\omega}}$ is assumed to be known. Note that Eq. (7) expresses the gradient of agent i 's ql- k Q -value function with respect to $\bar{\omega}$ in a form a Bellman equation that only depends on agent i 's ql- k Q -value function, the gradients of agents' reward functions with respect to $\bar{\omega}$, agents $-i$'s ql- $(k-1)$ policies and the corresponding policy gradients, which we all have access to. Therefore, we can compute $\frac{\partial Q_{\bar{\omega}}^{i,k}}{\partial \bar{\omega}}$ via value iteration algorithm and then $\frac{\partial \pi_{\bar{\omega}}^{i,k}}{\partial \bar{\omega}}$ can also be computed accordingly by differentiating Eq. (3).

Algorithm 1: Cognition-Aware Multi-Agent Inverse Reinforcement Learning algorithm

- 1 **Input:** A demonstration set \mathcal{D} and learning rate η
 - 2 **Output:** Learned parameters $\bar{\omega}$.
 - 3 Initialize $\bar{\omega}$.
 - 4 **while not converged do**
 - 5 Compute $\pi_{\bar{\omega}}^{i,k}$ and $\frac{\partial \pi_{\bar{\omega}}^{i,k}}{\partial \bar{\omega}}$ for each i and k ;
 - 6 Initialize $\nabla_{\bar{\omega}}$;
 - 7 **for** $\xi \in \mathcal{D}$ **do**
 - 8 Estimate agents' latent intelligence levels in ξ via Eq. (5);
 - 9 Compute gradient of the expected log-likelihood of the demonstrations following:

$$\nabla_{\bar{\omega}} + = \frac{\partial \left(\sum_{\bar{k} \in \mathbb{K}^n} \log \left(\mathbb{P}(\xi | \bar{\omega}, \bar{k}) \right) \mathbb{P}(\bar{k} | \xi, \bar{\omega}) \right)}{\partial \bar{\omega}};$$
 - 10 **end for**
 - 11 Update the parameters following: $\bar{\omega} = \bar{\omega} + \eta \nabla_{\bar{\omega}}$;
 - 12 **end while**
 - 13 **Return:** $\bar{\omega}$
-

Policy gradient. In the previous part, we show that with each agent's $\frac{\partial \pi_{\bar{\omega}}^{i,k-1}}{\partial \bar{\omega}}$ known, then we can compute $\frac{\partial \pi_{\bar{\omega}}^{j,k}}{\partial \bar{\omega}}$ for each agent $j \in \mathcal{P}$. Therefore, starting from $\frac{\partial \pi_{\bar{\omega}}^{i,0}}{\partial \bar{\omega}}$ of each agent, we can compute $\frac{\partial \pi_{\bar{\omega}}^{i,k}}{\partial \bar{\omega}}$ through a sequential and iterative process similar to the procedure for computing $\pi_{\bar{\omega}}^{i,k}$ desired be Section IV. Since ql-0 agents (agents with the lowest intelligence level) are normally represented as reflexive agents who do not consider sophisticated interactions with others [32], $\frac{\partial \pi_{\bar{\omega}}^{i,0}}{\partial \bar{\omega}}$ can be computed straightforwardly based on a particular instantiation. The procedure for computing $\frac{\partial \pi_{\bar{\omega}}^{i,k}}{\partial \bar{\omega}}$

is illustrated in lower-left block in Fig. 2.

With $\frac{\partial \pi^{i,k}}{\partial \bar{\omega}}$ computed, the gradient $\frac{\partial \mathbb{P}(k^i|\xi, \bar{\omega})}{\partial \bar{\omega}}$ can be obtained by differentiating Eq. (5), which yields a recursive format from time 0 to time t and can be easily computed with initialization $\frac{\partial \mathbb{P}(k^i|\xi_0, \bar{\omega})}{\partial \bar{\omega}} = \mathbf{0}$.

Learning algorithm. With both $\frac{\partial \pi^{i,k}}{\partial \bar{\omega}}$ and $\frac{\partial \mathbb{P}(k^i|\xi, \bar{\omega})}{\partial \bar{\omega}}$ available, the gradient of the learning objective function Eq. (4) with respect to $\bar{\omega}$ can be computed straightforwardly and used to update the estimate of $\bar{\omega}$. In practice, we also regularize the reward parameters during learning. The overview of learning curriculum is illustrated in Fig. 2 and our Cognition-Aware Multi-Agent Inverse Reinforcement Learning algorithm is summarized in Algorithm 1.

VI. EXPERIMENT DESIGN

In this section, we design two experiments to validate our hypothesis: *reasoning about humans' latent intelligence levels adds flexibility to inverse learning algorithm and helps learn reward functions that explain and reconstruct human driving behaviors better.*

A. Environments

Although our formulation and approach generally apply to settings with multiple agents ($n \geq 2$), here we focus on the cases where the ego agent interacts with another opponent agent in order to compare with baselines that are tailored for pair-wise relationships.

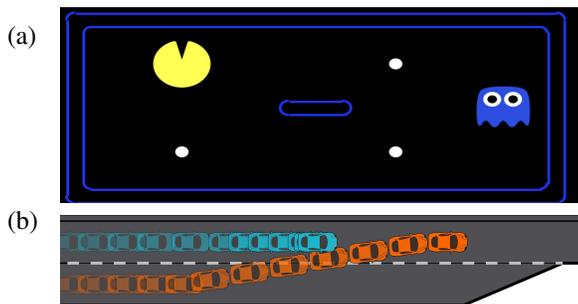


Fig. 3. Simulated environments. (a): Pac-man aims to collect all the dots and avoid being eaten by the ghost. (b): The lower-lane car wants to effectively execute the merging while remaining safe.

Zero-sum game. We consider a 2D maze game, Pac-Man (as shown in Fig. 3 (a)), in which the pac-man aims to collect all the pac-dots in the shortest time possible and the ghost aims to prevent the pac-man from doing so. The agents share the same action set that contains five actions: moving up, down, left, right, or stay. The state of the game is defined as $s = [x^1, y^1, x^2, y^2, \bar{\mathbb{I}}]$, where x (y) is the longitudinal (lateral) position, the superscript 1 (2) denotes the pac-man (ghost), and $\bar{\mathbb{I}} = [\mathbb{I}^1, \dots, \mathbb{I}^{n_d}]$ is a vector of Booleans that track the validity of the n_d pac-dots in the game.

General-sum game. We use a driving environment that consists of two cars on a two-lane road. As shown in Fig. 3 (b), the orange car has to merge to its adjacent (upper) lane which is occupied by the blue car. Sophisticated and interactive behaviors are likely to happen under such a scenario [33]. The dynamics of the driving scenario are represented as $[\dot{x}^1 \ \dot{y}^1 \ \dot{x}^2 \ \dot{v}^1 \ \dot{v}^2] = [v^1 \ w^1 \ v^2 \ a^1 \ a^2]$, where

x (y) is the longitudinal (lateral) position, v (w) is the longitudinal (lateral) speed, a is the longitudinal acceleration, and the superscript 1 (2) denotes the lower-lane (upper-lane) car. The sampling period is $\Delta t = 0.5[s]$. We use a state grid of size $50 \times 4 \times 50 \times 5 \times 5$ to represent the discrete states of the game in a similar way as in [18].

B. Manipulated Variables

The manipulated variable is the choice of different inverse reinforcement learning algorithms for learning agents' reward functions from demonstrations. In addition to our proposed approach, we consider two baseline algorithms that are both extensions of MaxEnt IRL in two-agent settings. In contrast to our approach, they both assign agents' roles/strategies *a priori* during learning.

PR2-MaxEnt IRL. Probabilistic recursive reasoning (PR2) framework [22] is a particular instantiation of the quantal level- k model. However, PR2 assumes all agents perform depth-1 recursive reasoning. As the first baseline, we use the PR2 framework to model interactions and plug the PR2- Q algorithm into our learning framework: we assume that all humans in a demonstration perform depth-1 recursive reasoning when interacting with others and we use the PR2- Q algorithm to generate humans' policies conditioned on the roll-out reward functions. Then, we use the generated policies to compute the gradient of the log-likelihood of demonstrations (Eq. (1)) without reasoning about humans' latent cognitive states since the PR2 framework assigns humans' intelligence levels *a priori*.

MaxEnt IRL with leader-follower (LF) roles. Our second baseline (LF-MaxEnt IRL) is the MaxEnt IRL algorithm used in [10], [12], [13], [17], which pre-assigns the leader and the follower during learning (refer to Section II for more details).

C. Dependent Measures

In order to validate our proposed approach and compare it against the baselines, we conducted two experiments. In the first experiment, we utilize demonstrations from synthetic agents in the two environments, then we can quantify the performance of IRL algorithms in terms of: (a) the correlations between the ground-truth reward functions (handcrafted) and the learned reward functions; (b) the log-likelihood of the demonstrations induced by an IRL algorithm. In the second experiment, we further evaluate the effectiveness different IRL algorithms in the driving domain with real traffic data by measuring: (a) the log-likelihood of the demonstrations induced by an IRL algorithm; (b) the similarity between the ground-truth interaction trajectories in test set and the trajectories generated using the reward functions learned by an IRL algorithm and its interaction model.

D. Feature Selection

To avoid confounding variables, we control the features associated with agents' reward functions in both environments. More specifically, we let agents optimize over a linear combination of common features in our approach and the baselines. The feature weights are our reward parameters $\bar{\omega}$.

Pac-Man environment. Pac-man considers the following features: 1) *Safety*: this feature represents the pac-man’s preference to avoid the ghost and is defined as $f_{1,p} = \exp(d)$, where d denotes the $L-1$ distance between the pac-man and the ghost. 2) *Proximity*: this feature represents how the pac-man wants to reach the closest pac-dot and is defined as $f_{2,p} = \frac{1}{l_{d,min}}$, where $l_{d,min}$ is the distance from the pac-man to the closet pac-dot. 3) *Winning*: this feature represents how the pac-man wants to eat more pac-dots for winning the game and is defined as $f_{3,p} = \frac{1}{n_{ud}}$, where n_{ud} denotes the number of uncollected pac-dots. The ghost considers similar features as the pac-man does but with modifications: $f_{1,g} = \frac{1}{d}$, $f_{2,g} = f_{2,p}$, and $f_{3,g} = \frac{1}{f_{3,p}}$.

Driving environment. For the driving domain, we select features based on previous works [13], [33]: 1) *Progress*: this feature represents humans’ preference to drive faster and merge to the desired lane. It is defined as $f_1 = (\frac{v}{v_{max}})^2 + (\frac{y_{eff}}{|y_l - y_0|})^2$, where v_{max} denotes the maximum allowable speed, y_0 denotes the initial lateral position of the car, y_l denotes the lateral coordinate of the target lane, and y_{eff} denotes the effective transnational distance that the car has executed towards the target lane ($\max(\frac{y_{eff}}{|y_l - y_0|}) = 1$ and $\min(\frac{y_{eff}}{|y_l - y_0|}) = 0$). 2) *Comfort*: this feature represents humans’ desire to operate smoothly and is defined as $f_2 = (\frac{a}{a_{max}})^2$. 3) *Safety*: this feature represents how humans want to avoid collisions, and is defined as $f_3 = \sum_{i=1}^{n_o} (\frac{\min(d_i, d_{safe})}{d_{safe}})^2$, where d_{safe} is a pre-defined safety distance and d_i is the distance between the ego agent and the i -th obstacle.

VII. RESULTS AND ANALYSIS

A. Performance with Demonstrations from Synthetic Agents

Synthetic agents. In both environments, we model the synthetic agents as ql- k agents with various $k \in \mathbb{K}$. Recall that policies of ql-0 agents are required to initiate the recursive reasoning process Eq. (2). In line with previous works [20], [23], [34], we let ql-0 agents be reflective (non-strategic) agents who do not explicitly take into account their opponents’ possible responses but rather maximize their immediate rewards by treating other agents as stationary objects. The ground-truth reward functions of synthetic agents in each environment are manually tuned to achieve reasonable behaviors. For each environment, we generate 30 interactions with random initial states and ground-truth agent types (k).

Learning performance. In Fig. 4, we show the histories of the log-likelihood of the demonstration set (learning objective function) using our algorithm in each environment (blue line). The black lines denote the ground truth log-likelihood of the demonstration set evaluated using the models of synthetic agents. It can be observed that the log-likelihood of the demonstration set induced by the learned reward functions approaches to the ground-truth value as the learning algorithm converges. In Fig. 5, we compare the learned reward parameters and the ground-truth reward parameters using Pearson’s correlation coefficient (PCC) and Spearman’s rank correlation coefficient (SCC). In general,

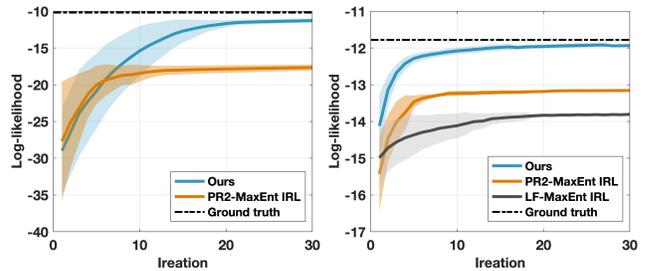


Fig. 4. Histories of the log-likelihood of the demonstration set during learning (solid line represents the mean and the light shaded area represents the 95% confidence tube of the data). Left: Pac-Man. Right: Driving.

a higher PCC indicates a higher linear correlation, and a higher SCC represents a stronger monotonic relationship. It can be observed that our approach is able to recover reward parameters with a high linear correlation and a strong monotonic relationship to the ground-truth ones.

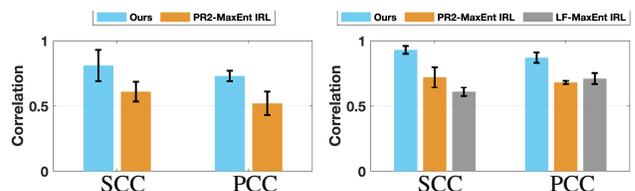


Fig. 5. Statistical correlations between the learned reward parameters and the ground-truth ones. Left: Pac-Man. Right: Driving.

Remark. In Fig. 4 and Fig. 5, we also plot the results with the baselines. Such a comparison with the baselines mainly serves as a sanity check: with synthetic agents, our approach is supposed to work better than the baselines since the interaction models of synthetic agents align with our approach to modeling interactions during learning. However, this experiment illustrates the effectiveness of our algorithm in both adversarial and cooperative settings if agents were indeed following our Theory-of-Mind model for making decisions. In the following section, we utilize real human driving data to test the performance of our approach and make a detailed comparison against the baselines.

B. Learning Driver Reward Functions from Traffic Data

Traffic data. We extract real driving data in a forced merging scenario (DRCHNMergingZS) from INTERACTION dataset [35]. We extract 30 interactions (sampling time: 0.5[s]) as the training set and another 15 interactions as the test set.

Learning performance. In the left plot of Fig. 6, we show the histories of the log-likelihood of the demonstrations during learning. It can be observed that our approach can better explain human driving behaviors, yielding a higher likelihood for the demonstration set based on the learned reward functions and the exploited interaction model. More importantly, our approach performs better at very early iterations, which indicates that, by reasoning about human drivers’ latent cognitive state (intelligence level), the learning algorithm is offered more flexibility to explain human behaviors.

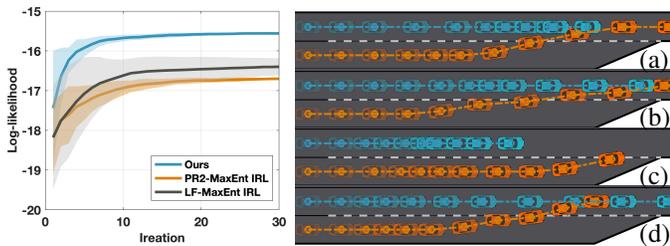


Fig. 6. Left: Histories of the log-likelihood of the demonstration set during learning. Right: An example of test interaction trajectory and reconstructed interaction trajectories (a: ground truth; b: Ours; c: PR2-MaxEnt IRL; d: LF-MaxEnt IRL).

Analysis of the learned driving preferences. In Fig. 7, we show the convergence of reward parameters (feature weights) in 15 trials with random initialization. It can be observed that our approach learns that the lower-lane car (orange car) balances progress, comfort, and safety (Fig. 7(a)), while the upper-lane car (blue car) emphasizes more on progress and comfort (Fig. 7(b)). Such results align well with our experience and intuition that drivers with right-of-way tend to shift safety responsibilities to others without right-of-way. With PR2-MaxEnt IRL (Fig. 7(c-d)), due to the homogeneous assumption on the interaction model, the learned weights for both cars are biased towards safety and comfort. With LF-MaxEnt IRL, the weights of the lower-lane car are biased towards safety and comfort (Fig. 7(e)), and the weights of the upper-lane car converge to two clusters, with one emphasizing more on safety and the other more on progress (Fig. 7(f)). This is because that LF-MaxEnt IRL pre-assigns the roles of humans, thus it is unable to capture the structural bias in human mind, leading to weights with a large variance.

TABLE I
SIMILARITY SCORES OF DIFFERENT ALGORITHMS

| Algo. | Ours | PR2-MaxEnt IRL | LF-MaxEnt IRL |
|-------------|--------------|----------------|---------------|
| Traj. Score | 10.86 | 17.38 | 24.02 |
| Deci. Score | 0.93 | 0.73 | 0.6 |

Analysis of the re-generated interaction trajectories. Due to the unknown ground-truth reward functions, we investigate whether our approach and the baselines can reproduce interactions in the test set using the learned reward functions and their interaction models. In the right plot of Fig. 6, we show an example of the ground-truth interaction trajectory (a) in the test set, and the re-generated interaction trajectories by our approach (b) and the baselines (c-d). The test interaction example is intentionally selected to be complex: both drivers accelerate initially, then the lower-lane driver decelerates slightly but accelerates to merge after observing the deceleration from the upper-lane driver. Our approach is able to generate a seamless interaction that reproduces the behaviors in the test interaction. On the contrary, the baselines are unable to reproduce a similar interaction due to their biased assumptions. Specifically, PR2-MaxEnt IRL assumes that all agents in the game perform L-1 recursive reasoning, thus a dead-lock behavior emerged in the generated interaction: both agents tend to yield initially, then the lower-lane driver initiates the merge forced by the approaching dead-end.

The leading-following interaction model exploited by LF-MaxEnt IRL assumes one agent aims to compute the best response to the known trajectory of the other agent, thus in the interaction reproduced by LF-MaxEnt IRL, both drivers try to accommodate their opponents' ground-truth trajectories, leading to an incorrect interaction (lower-lane driver yields to the upper-lane driver). In Table I, we quantitatively compare the similarity between the re-generated interaction trajectories and the ground-truth ones, at both the trajectory level and the decision level. At the trajectory level, we define similarity score as the (averaged) Euclidean distance between the reproduced trajectories and test trajectories. At the decision level, we define the similarity score as the accuracy of the agents' high-level decisions (yield or not) in the reproduced trajectories with the high-level decisions in the test trajectories as ground-truth decisions. It can be observed from Table I that our approach performs significantly better than the baselines.

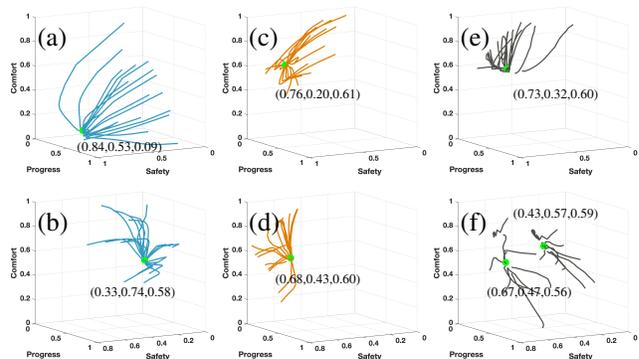


Fig. 7. Histories of feature weights during learning (each line represents a data trail that shows the trajectory of weights during one learning trial). Top row: orange car. Bottom row: blue car. (a-b): Ours. (c-d): PR2-MaxEnt IRL. (e-f): LF-MaxEnt IRL. The green points denote the converged weights in each cluster, and the texts around the green points show the values of the converged weights (safety, progress, comfort).

VIII. DISCUSSION

Summary. In this work, we advocated that humans have different levels of sophistication in reasoning about others' behaviors during interactions, and such an aspect should be accounted for during reward function learning. We exploited insights from Theory-of-Mind and proposed a new Maximum Entropy-based multi-agent Inverse Reinforcement Learning framework that reasons about humans' latent intelligence levels during learning. We validated our approach in both zero-sum and general-sum games with ground truth humans, then illustrated a practical application of our approach to learning human drivers' reward functions from real driving data. We compared our approach with the baseline algorithms and provided a detailed experiment analysis. We found that reasoning about humans' levels of intelligence provides more flexibility during learning and helps explain and reconstruct human driving behaviors better.

Limitations and future works. We view our work as a first step into incorporating Theory-of-Mind based bounded intelligence reasoning into multi-agent IRL. One of our approach's limitations is the ability to treat continuous states

and actions. Fortunately, the procedure for computing $\pi^{i,k}$ can be realized in an iterative deep Q -learning fashion [36] and plugged into Algorithm 1 seamlessly. Another limitation is related to the assumption about humans' constant intelligence levels during interactions. Such an assumption is commonly used in one-shot games. However, by inspecting traffic data, we noticed that there might be a short period at the beginning of an interaction during which the agents compete for higher levels of intelligence in order to dominate the interaction. We envision that incorporating such a dynamic attribute of cognitive states into reward learning could potentially gain more flexibility during learning.

REFERENCES

- [1] S. Russell, "Learning agents for uncertain environments," in *Proceedings of the eleventh annual conference on Computational learning theory*, 1998, pp. 101–103.
- [2] A. Y. Ng, S. J. Russell, et al., "Algorithms for inverse reinforcement learning," in *Icml*, vol. 1, 2000, pp. 663–670.
- [3] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *Proceedings of the twenty-first international conference on Machine learning*, 2004, p. 1.
- [4] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning," in *Aaai*, vol. 8. Chicago, IL, USA, 2008, pp. 1433–1438.
- [5] C. Finn, P. Christiano, P. Abbeel, and S. Levine, "A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models," *arXiv preprint arXiv:1611.03852*, 2016.
- [6] A. I. Goldman et al., "Theory of mind," *The Oxford handbook of philosophy of cognitive science*, vol. 1, 2012.
- [7] L. Yu, J. Song, and S. Ermon, "Multi-agent adversarial inverse reinforcement learning," in *International Conference on Machine Learning*. PMLR, 2019, pp. 7194–7201.
- [8] N. Gruver, J. Song, M. J. Kochenderfer, and S. Ermon, *Multi-Agent Adversarial Inverse Reinforcement Learning with Latent Variables*. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, 2020, p. 1855–1857.
- [9] G. Coricelli and R. Nagel, "Neural correlates of depth of strategic reasoning in medial prefrontal cortex," *Proceedings of the National Academy of Sciences*, vol. 106, no. 23, pp. 9163–9168, 2009.
- [10] D. Sadigh, A. D. Dragan, S. Sastry, and S. A. Seshia, "Active preference-based learning of reward functions," in *Robotics: Science and Systems*, 2017.
- [11] E. Biyik and D. Sadigh, "Batch active preference-based learning of reward functions," in *Conference on Robot Learning*. PMLR, 2018, pp. 519–528.
- [12] L. Sun, W. Zhan, and M. Tomizuka, "Probabilistic Prediction of Interactive Driving Behavior via Hierarchical Inverse Reinforcement Learning," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, Nov. 2018, pp. 2111–2117.
- [13] L. Sun, W. Zhan, M. Tomizuka, and A. D. Dragan, "Courteous autonomous cars," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 663–670.
- [14] D. Sadigh, S. S. Sastry, S. A. Seshia, and A. Dragan, "Information gathering actions over human internal state," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 66–73.
- [15] S. V. Albrecht, S. Liemhetcharat, and P. Stone, "Special issue on multiagent interaction without prior coordination: Guest editorial," *Autonomous Agents and Multi-Agent Systems*, vol. 31, no. 4, pp. 765–766, 2017.
- [16] D. Fridovich-Keil, A. Bajcsy, J. F. Fisac, S. L. Herbert, S. Wang, A. D. Dragan, and C. J. Tomlin, "Confidence-aware motion prediction for real-time collision avoidance," *The International Journal of Robotics Research*, vol. 39, no. 2-3, pp. 250–265, 2020.
- [17] D. Sadigh, S. Sastry, S. A. Seshia, and A. D. Dragan, "Planning for autonomous cars that leverage effects on human actions," in *Robotics: Science and Systems*, vol. 2. Ann Arbor, MI, USA, 2016.
- [18] J. F. Fisac, E. Bronstein, E. Stefansson, D. Sadigh, S. S. Sastry, and A. D. Dragan, "Hierarchical game-theoretic planning for autonomous vehicles," in *2019 International Conference on Robotics and Automation (ICRA)*, May 2019, pp. 9590–9596.
- [19] C. Baker, R. Saxe, and J. Tenenbaum, "Bayesian theory of mind: Modeling joint belief-desire attribution," in *Proceedings of the annual meeting of the cognitive science society*, vol. 33, no. 33, 2011.
- [20] N. Li, D. W. Oyler, M. Zhang, Y. Yildiz, I. Kolmanovsky, and A. R. Girard, "Game theoretic modeling of driver and vehicle interactions for verification and validation of autonomous vehicle control systems," *IEEE Transactions on control systems technology*, vol. 26, no. 5, pp. 1782–1797, 2017.
- [21] H. de Weerd, R. Verbrugge, and B. Verheij, "Negotiating with other minds: the role of recursive theory of mind in negotiation with incomplete information," *Autonomous Agents and Multi-Agent Systems*, vol. 31, no. 2, pp. 250–287, 2017.
- [22] Y. Wen, Y. Yang, R. Luo, J. Wang, and W. Pan, "Probabilistic recursive reasoning for multi-agent reinforcement learning," *arXiv preprint arXiv:1901.09207*, 2019.
- [23] R. Tian, N. Li, I. Kolmanovsky, Y. Yildiz, and A. R. Girard, "Game-theoretic modeling of traffic in unsignalized intersection network for autonomous vehicle control verification and validation," *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [24] L. Sun, Z. Wu, H. Ma, and M. Tomizuka, "Expressing diverse human driving behavior with probabilistic rewards and online inference," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 2020–2026.
- [25] Z. Wu, L. Sun, W. Zhan, C. Yang, and M. Tomizuka, "Efficient sampling-based maximum entropy inverse reinforcement learning with application to autonomous driving," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5355–5362, 2020.
- [26] F. Memarian, Z. Xu, B. Wu, M. Wen, and U. Topcu, "Active task-inference-guided deep inverse reinforcement learning," in *2020 59th IEEE Conference on Decision and Control (CDC)*. IEEE, 2020, pp. 1932–1938.
- [27] J. K. Goeree and C. A. Holt, "Ten little treasures of game theory and ten intuitive contradictions," *American Economic Review*, vol. 91, no. 5, pp. 1402–1422, 2001.
- [28] V. P. Crawford and N. Iriberry, "Level-k auctions: Can a nonequilibrium model of strategic thinking explain the winner's curse and overbidding in private-value auctions?" *Econometrica*, vol. 75, no. 6, pp. 1721–1770, 2007.
- [29] R. Tian, L. Sun, and M. Tomizuka, "Bounded risk-sensitive markov games: Forward policy design and inverse reward learning with iterative reasoning and cumulative prospect theory," in *AAAI Conference on Artificial Intelligence*, 2021.
- [30] D. O. Stahl II and P. W. Wilson, "Experimental evidence on players' models of other players," *Journal of economic behavior & organization*, vol. 25, no. 3, pp. 309–327, 1994.
- [31] J. R. Wright and K. Leyton-Brown, "Predicting human behavior in unrepeated, simultaneous-move games," *Games and Economic Behavior*, vol. 106, pp. 16 – 37, 2017. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0899825617301574>
- [32] —, "Level-0 meta-models for predicting human behavior in games," in *Proceedings of the fifteenth ACM conference on Economics and computation*, 2014, pp. 857–874.
- [33] R. Choudhury, G. Swamy, D. Hadfield-Menell, and A. D. Dragan, "On the utility of model learning in hri," in *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 2019, pp. 317–325.
- [34] M. Bouton, A. Nakhai, D. Isele, K. Fujimura, and M. J. Kochenderfer, "Reinforcement learning with iterative reasoning for merging in dense traffic," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2020, pp. 1–6.
- [35] W. Zhan, L. Sun, D. Wang, H. Shi, A. Clausse, M. Naumann, J. Kümmerle, H. Königshof, C. Stiller, A. de La Fortelle, and M. Tomizuka, "INTERACTION Dataset: An INTERNATIONAL, Adversarial and Cooperative MOTION Dataset in Interactive Driving Scenarios with Semantic Maps," *arXiv:1910.03088 [cs, eess]*, Sept. 2019.
- [36] M. Bouton, A. Nakhai, D. Isele, K. Fujimura, and M. J. Kochenderfer, "Reinforcement learning with iterative reasoning for merging in dense traffic," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, 2020, pp. 1–6.