

# INVERSE REINFORCEMENT LEARNING OF INTERACTIVE SCENARIOS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This paper studies the problem where a learner aims to learn the reward function of an expert and a policy to interact with the expert from interactions with the expert. We formulate the problem as a stochastic bi-level optimization problem where the lower level learns a reward function that explains the behaviors of the expert, and the upper level learns a policy to interact with the expert. We develop a double-loop algorithm, General Scenario Interactive Inverse Reinforcement Learning (GSIIRL), which solves the lower-level optimization problem in the inner loop and the upper-level optimization problem in the outer loop. We formally guarantee that GSIIRL converges at the rate of  $\mathcal{O}(1/\sqrt{K})$  and empirically validate our algorithm through simulations.

## 1 INTRODUCTION

Inverse reinforcement learning (IRL) aims to recover a reward function and a corresponding policy that best explain an expert’s behavior given the expert’s demonstration trajectories. IRL has been applied across diverse domains, including robotics (Ziebart et al., 2008; Okal & Arras, 2016), cybersecurity (Zhang et al., 2019; Elnaggar & Bezzo, 2018), and biology (Hirakawa et al., 2018; Ashwood et al., 2022). In standard IRL, the expert’s demonstration trajectories are assumed to be fixed and unaffected by the learner, as the learner passively observes the expert and learns from these demonstrations without influencing the expert’s behavior. However, recent applications (Palaniappan et al., 2017; Büning et al., 2022; Kamalaruban et al., 2019) have motivated an interactive form of IRL where the learner actively interacts with the expert to infer the reward function and policy through these interactions. For example, a robot needs to infer a human’s intended destination in a maze and then open the appropriate doors to help the human reach that goal (Büning et al., 2022). In this scenario, the learner (robot)’s actions (e.g., opening different doors) directly influence the expert (human)’s trajectory, violating the standard IRL assumption of passive observation and rendering conventional IRL methods inapplicable.

To bridge this gap, this paper studies interactive IRL, where a learner interacts with an expert to learn a reward function that explains the expert’s behavior (i.e., such that the expert’s trajectories are optimal under the learned reward) and a policy for effectively interacting with the expert. In this interactive setting, the learner is not merely an observer but an active participant: by interacting with the expert, the learner can influence the expert’s trajectories. Consequently, if the learner’s policy changes, the expert’s trajectories may change accordingly. Consider a scenario where a human (expert) and a ground mobile robot (learner) share a 2D environment while navigating to their respective destinations. The human has the right-of-way, so the robot must either slow down or detour to allow the human to cross first. At the same time, the robot is trying to reach its own destination as soon as possible. Thus, the robot has a human-independent objective (e.g., reaching its destination) as well as a human-dependent objective (e.g., bypassing the human). Accordingly, the robot’s reward function can be decomposed into a human-independent part and a human-dependent part. For example, in (El-Shamouty et al., 2020), the reward function is defined by two terms: the robot’s distance to its goal and whether a collision with the human occurs.

Motivated by this example, this paper assumes that the learner’s reward function can be decomposed into an expert-independent component and an expert-dependent component. In our interactive IRL framework, the relationship between the learner’s and the expert’s reward functions is more general than those in fully cooperative interactive IRL (Palaniappan et al., 2017; Büning et al., 2022; Ka-

054 malaruban et al., 2019; Hadfield-Menell et al., 2016) and fully competitive interactive IRL (Zhang  
055 et al., 2019; Wang & Klabjan, 2018) settings. In the fully cooperative case, the learner and the ex-  
056 pert share an identical reward function, whereas in the fully competitive case, the learner’s reward  
057 function is the negation of the expert’s reward function.

058 The interactive IRL problem of interest is closely related to RL (Lowe et al., 2017; Yu et al., 2022;  
059 Kuba et al., 2021), and classic IRL (Lin et al., 2019; Yu et al., 2019; Liu & Zhu, 2022). However,  
060 standard RL algorithms are not applicable to our setting because they require access to state-action-  
061 reward tuples of the expert, which are unavailable in interactive IRL. Similarly, classical IRL al-  
062 gorithms require passive observation, whereas in our setting, the learner actively interacts with the  
063 expert. These limitations also apply to multi-agent RL (MARL) and multi-agent IRL (MA-IRL),  
064 respectively. Detailed comparisons will be further elaborated in Section 3.

065 **Contribution.** This paper makes four main contributions. First, we study an interactive IRL prob-  
066 lem, where the learner learns a reward function to explain the expert’s behaviors and a policy to  
067 interact with the expert. We formulate this as a stochastic bi-level optimization problem in which  
068 the lower level learns a reward function to explain the expert’s behavior and the upper level learns  
069 a policy to interact with the expert. Second, we develop a double-loop algorithm, General Scenario  
070 Interactive Inverse Reinforcement Learning (GSIIRL), to solve the bi-level optimization problem,  
071 with an outer loop handling the upper-level optimization problem and an inner loop handling the  
072 lower-level optimization problem. A key challenge in this approach is computing the hypergradient  
073 (the gradient of the upper-level objective function), since it involves an intractable Hessian. We  
074 address this by using simultaneous perturbation stochastic approximation (SPSA) (Spall, 1992) to  
075 approximate the hypergradient, reducing the computational complexity per outer-loop iteration from  
076  $\mathcal{O}(T^2)$  to  $\mathcal{O}(T)$ , where  $T$  is the maximum trajectory length. Third, we theoretically guarantee that  
077 the algorithm converges at a rate of  $\mathcal{O}(1/\sqrt{K})$ , where  $K$  is the number of outer-loop iterations.  
078 Fourth, we validate our algorithm through four experiments and show that GSIIRL performs com-  
079 parably to benchmark methods (MARL and MA-IRL) that require additional information (either the  
080 expert’s reward function or learner’s optimal demonstrations).

## 081 2 RELATED WORK

082 **IRL and multiagent IRL (MA-IRL).** As (Ng & Russell, 2000) mentioned, IRL faces the challenge  
083 that the demonstrated trajectories can be explained by multiple reward functions. Several approaches  
084 have been proposed to address this challenge. The current state-of-the-art IRL methods include max-  
085 imum margin IRL (Ng & Russell, 2000; Abbeel & Ng, 2004), maximum entropy IRL (Ziebart et al.,  
086 2008; 2010), maximum likelihood IRL (Zeng et al., 2022; 2023), and Bayesian IRL (Ramachandran  
087 & Amir, 2007; Choi & Kim, 2012). However, these methods all focus on single-expert scenarios.  
088 MA-IRL extends IRL to settings with multiple experts, where one or more learners recover the re-  
089 ward functions of multiple experts from their demonstrations (Lin et al., 2019; Yu et al., 2019; Liu  
090 & Zhu, 2022). These MA-IRL approaches still assume that the learners are isolated from the experts  
091 (i.e., no direct interaction between them) and thus cannot address interactive IRL scenarios.

092 **Fully cooperative and fully competitive interactive IRL.** Prior work on interactive IRL has pri-  
093 marily explored two cases: fully cooperative and fully competitive scenarios. In the fully cooperative  
094 setting, the learner and expert share the same reward function, and their interaction is formulated  
095 as a cooperative game (Palaniappan et al., 2017; Büning et al., 2022; Kamalaruban et al., 2019;  
096 Hadfield-Menell et al., 2016). In this case, the learner’s objective is to recover the expert’s reward  
097 function. In contrast, the fully competitive setting assumes the learner’s reward is the negative of the  
098 expert’s reward, leading to a zero-sum game (Zhang et al., 2019; Wang & Klabjan, 2018). Here, the  
099 learner aims to learn an objective that directly opposes the expert’s reward. Both of these approaches  
100 impose strong assumptions on the relationship between the expert’s and learner’s reward functions,  
101 specifically assuming they are either identical or exact opposites. In this paper, the relationship  
102 between the reward functions of the learner and the expert is arbitrary.

103 **Bi-level optimization.** Bi-level optimization has been applied to many machine learning prob-  
104 lems, including meta-learning (Lee et al., 2019; Xu & Zhu, 2022), hyperparameter optimization  
105 (Pedregosa, 2016; Xu & Zhu, 2023), and IRL (Liu & Zhu, 2022; 2024). A classic approach to such  
106 problems is the descent method (Kolstad & Lasdon, 1990; Xu & Zhu, 2023). This method typically  
107 requires computing the second-order Hessian of the lower-level objective function, an operation that

is prohibitively expensive to compute in our setting. A common strategy to avoid explicitly computing the Hessian is to approximate it using finite differences (Strikwerda, 2004). In this paper, we further reduce the computational burden by adopting the SPSA method, which perturbs all dimensions of the decision variable simultaneously and therefore requires only two objective function evaluations per iteration.

### 3 MODEL AND PROBLEM STATEMENT

In this section, we introduce the interactive IRL problem.

**Markov Game model.** We model the interactions between a learner and an expert as a finite horizon Markov Game (MG)  $(S, A, P, T, r_{ld} + r_{li}, r_e, \gamma)$ . The elements of the MG are defined as follows:

- $S \triangleq S_l \times S_e$  is the continuous state space where  $S_l$  and  $S_e$  are the state spaces of the learner and the expert, respectively. We denote  $s = (s_l, s_e) \in S$ , with  $s_l \in S_l$  and  $s_e \in S_e$ .
- $A \triangleq A_l \times A_e$  is the joint action space, where  $A_l$  and  $A_e$  are the continuous action spaces of the learner and the expert, respectively. We denote  $a = (a_l, a_e) \in A$ , with  $a_l \in A_l$  and  $a_e \in A_e$ .
- $P(s'|s, a)$  is the transition probability density for moving from state  $s$  to  $s'$  by taking joint action  $a$ .
- $T$  is the finite time horizon.
- $r_{ld} + r_{li}$  is the reward function of the learner, consisting of two parts: the expert-dependent reward function  $r_{ld}$  and the expert-independent reward function  $r_{li}$ . The function  $r_{ld}$  maps full state-action pairs  $(s, a)$  to bounded rewards, and  $r_{li}$  maps learner-only state-action pairs  $(s_l, a_l)$  to bounded rewards.
- $r_e$  is the expert's reward function, mapping state-action pairs  $(s, a)$  to bounded rewards.
- $\gamma \in (0, 1]$  is the discount factor.

We define  $\pi_l(a_l|s)$  as the learner's policy and  $\pi_e(a_e|s)$  as the expert's policy. The joint policy is defined as  $\pi(a|s) \triangleq \pi_l(a_l|s) \times \pi_e(a_e|s)$ , which represents the probability density of the learner taking action  $a_l$  and the expert taking  $a_e$  at state  $s$ . When the joint policy  $\pi$  is executed, the MG generates a trajectory  $\zeta = s^0, a^0, s^1, a^1, \dots, s^{(T-1)}, a^{(T-1)}$ .

The learner's policy  $\pi_l$  aims to maximize the cumulative reward  $E^{\pi_l, \pi_e} [\sum_{t=0}^{T-1} \gamma^t (r_{ld}(s^t, a^t) + r_{li}(s_l^t, a_l^t))]$ . Analogously, the expert's policy  $\pi_e$  aims to maximize  $E^{\pi_l, \pi_e} [\sum_{t=0}^{T-1} \gamma^t r_e(s^t, a^t)]$ .

**Knowledge and goal of the learner.** The learner does not know the expert's reward function  $r_e$  but knows its own expert-independent reward function  $r_{li}$ . It can observe state-action pairs  $(s, a)$  and interact with the expert. Based on the trajectories collected during the interactions, the learner aims to recover the expert's reward function  $r_e$  and compute the optimal joint policy  $\pi^*$ .

**Benefits of learning reward functions.** As noted by the seminal work (Ng & Russell, 2000), "the reward function is the most succinct, robust, and transferable definition of the task". When the transition probabilities of a Markov Game substantially change, the policy needs to be retrained, but the underlying reward function can remain the same. Furthermore, learning a human's reward function can facilitate the robot's policy learning (Büning et al., 2022). If the ground-truth  $r_{ld}$  and  $r_e$  are sparse, learning dense reward functions to approximate them can accelerate the policy learning while leaving the optimal policy unchanged (Memarian et al., 2021). Given the importance of reward functions, we include learning  $r_e$  as a central goal for the learner.

We next discuss the distinctions between interactive IRL and related problems, including RL, IRL, fully cooperative interactive IRL, and fully competitive interactive IRL.

**Distinctions from RL.** In our interactive IRL setting, the learner's objective explicitly includes learning the expert's reward function  $r_e$  from observed expert trajectories. These trajectories contain implicit information about  $r_e$  because they are generated by the expert's policy  $\pi_e$ , which maximizes the expected cumulative reward  $E^{\pi_l, \pi_e} [\sum_{t=0}^{T-1} \gamma^t r_e(s^t, a^t)]$ . Model-free RL algorithms are not designed to learn reward functions. Model-based RL algorithms, on the other hand, can learn the expert's reward function given complete state-action-reward tuples. However, in our case, the learner only observes state-action pairs without access to the expert's reward signals, making model-based RL approaches also inapplicable. Similarly, MARL algorithms, as extensions of standard RL to multi-agent settings, inherit this limitation and thus cannot address interactive IRL problems.

**Distinctions from IRL.** In classical IRL, the learner passively observes expert demonstrations and has no influence over the expert’s behavior. The expert’s policy  $\pi_e'$  results from maximizing the expected cumulative reward  $E^{\pi_e'}[\sum_{t=0}^{T-1} \gamma^t r_e(s_e^t, a_e^t)]$ . Specifically, the expert’s action depends solely on its own state, and its trajectories are determined entirely by its policy and the environment dynamics. In contrast, interactive IRL is formulated within an MG framework, where the learner and expert interact with each other. In this setting, the expert’s actions depend on the joint state of both the expert and the learner. Thus, the learner’s policy influences its own state, which in turn affects the joint state, thereby affecting the expert’s actions and resulting trajectories. Here, the expert’s policy  $\pi_e$  maximizes the expected cumulative reward  $E^{\pi_l, \pi_e}[\sum_{t=0}^{T-1} \gamma^t r_e(s^t, a^t)]$ , explicitly incorporating the learner’s policy  $\pi_l$ . Due to these fundamental differences, classical IRL algorithms are unsuitable for interactive IRL problems, including ours, fully cooperative IRL (Palaniappan et al., 2017; Büning et al., 2022; Kamalaruban et al., 2019; Hadfield-Menell et al., 2016) and fully competitive IRL (Zhang et al., 2019; Wang & Klabjan, 2018). MA-IRL, as an extension of IRL to multi-agent settings, inherently shares this limitation and thus cannot address interactive IRL problems.

**Distinctions from fully cooperative and fully competitive interactive IRL.** Fully cooperative interactive IRL settings (Palaniappan et al., 2017; Büning et al., 2022; Kamalaruban et al., 2019; Hadfield-Menell et al., 2016) assume that the learner’s and expert’s reward functions are identical. Fully competitive interactive IRL settings (Zhang et al., 2019; Wang & Klabjan, 2018) assume that the learner’s reward function is the negative of the expert’s reward function. In our interactive IRL setting, the relationship between the learner’s and the expert’s reward functions can be arbitrary, which encompasses both identical and opposite cases.

#### 4 PROBLEM FORMULATION AND BI-LEVEL SETUP

In this section, we formulate the learning problem in Section 3 as a bi-level optimization problem where the lower-level optimization problem learns the expert reward function  $r_e$  and the upper-level optimization problem learns the optimal joint policy  $\pi^*$ .

Recall that  $r_e$  and  $r_{ld}$  are unknown to the learner. The learner uses a parametric function  $r_{\theta_e}$  to estimate  $r_e$  and uses  $r_{\theta_l}$  to estimate  $r_{ld}$  where  $\theta_e, \theta_l \in \Theta \triangleq \{\theta \mid \|\theta\|_2 \leq 1\}$ .

We define the policy  $\pi_{\theta_l, r_e}$  as the optimal joint policy when the learner’s reward function is  $r_{\theta_l} + r_{li}$  and the expert’s reward function is  $r_e$ . Analogously, given reward functions  $r_{\theta_l} + r_{li}$  and  $r_{\theta_e}$  for the learner and the expert respectively, the optimal joint policy is  $\pi_{\theta_l, \theta_e}$ .

**Learning expert’s reward function  $r_e$ .** The learner aims to infer the expert’s reward function according to the given estimate  $r_{\theta_l}$ . During interactions between the learner (using  $r_{\theta_l}$ ) and the expert (using the ground-truth reward function  $r_e$ ), the learner collects a set of trajectories  $D_{\theta_l, r_e} \triangleq \{\zeta_i\}_{i=1}^d$  generated by the policy  $\pi_{\theta_l, r_e}$ . Since the expert utilizes the ground-truth reward function  $r_e$ , the trajectories recorded from the interactions are the demonstrations of the expert. Given that the learner knows  $r_{\theta_l}$ , inferring  $r_e$  from collected demonstrations is a special case of the standard two-agent IRL problem (Lin et al., 2019; Yu et al., 2019; Liu & Zhu, 2022) with one of the reward functions known to the learner. This learning procedure, depicted in the upper block of Figure 1, can be formulated as the following maximum likelihood IRL (ML-IRL) problem:

$$\theta_e^*(\theta_l) = \arg \min_{\theta_e \in \Theta} L(\theta_l, \theta_e), \quad (1)$$

where the loss function is  $L(\theta_l, \theta_e) \triangleq -\sum_{i=1}^d \sum_{t=0}^{T-1} [\ln \pi_{\theta_l, \theta_e}(a^{it} | s^{it})] + \frac{\lambda}{2} \|\theta_e\|_2^2$ ,  $(a^{it}, s^{it}) \in \zeta_i \in D_{\theta_l, r_e}$ . This ML-IRL approach seeks a reward parameter  $\theta_e^*(\theta_l)$  such that the associated reward function  $r_{\theta_e^*(\theta_l)}$  maximizes the likelihood of generating the observed trajectories in  $D_{\theta_l, r_e}$ . With a

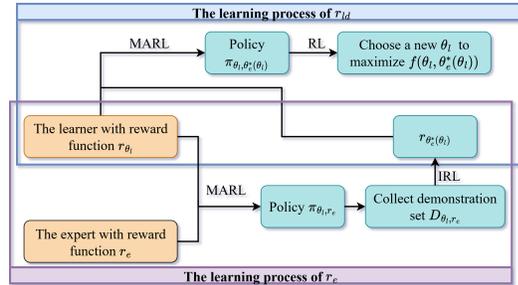


Figure 1: Flowchart of the overall learning process. The process of learning  $r_{ld}$  is included in the upper block, and the process of learning  $r_e$  is included in the lower block.

given  $r_{\theta_l}$ , the policy  $\pi_{\theta_l, \theta_e^*(\theta_l)}$  induced by the learned expert’s reward function  $r_{\theta_e^*(\theta_l)}$  is expected to produce the same trajectory from the policy  $\pi_{\theta_l, r_e}$ . Notice that the maximum likelihood estimation has been extensively applied in IRL (Ziebart et al., 2008; 2010). Additionally, to promote simpler reward structures, an  $L_2$  regularization term  $\frac{\lambda}{2} \|\theta_e\|_2^2$  with a small positive scalar  $\lambda$  is incorporated.

**Learning the joint policy  $\pi^*$  and the expert-dependent reward function  $r_{ld}$ .** Given a learner’s reward parameter  $\theta_l$ , the learner solves problem (1) to obtain the expert’s reward parameter  $\theta_e^*(\theta_l)$  and the corresponding policy  $\pi_{\theta_l, \theta_e^*(\theta_l)}$ . The learner’s objective is to learn a parameter  $\theta_l$  such that the policy  $\pi_{\theta_l, \theta_e^*(\theta_l)}$  maximizes the cumulative reward  $E^{\pi_{\theta_l, \theta_e^*(\theta_l)}} [\sum_{t=0}^{T-1} \gamma^t (r_{ld}(s^t, a^t) + r_{li}(s_t^t, a_t^t))]$ . This learning process is depicted in the lower block of Figure 1 and formulated as follows:

$$\arg \min_{\theta_l \in \Theta} f(\theta_l, \theta_e^*(\theta_l)) \triangleq -E^{\pi_{\theta_l, \theta_e^*(\theta_l)}} \left[ \sum_{t=0}^{T-1} \gamma^t (r_{ld}(s^t, a^t) + r_{li}(s_t^t, a_t^t)) \right] \quad (2)$$

To solve the problem (2), the learner needs to interact with the expert and receive the reward values of  $r_{ld}(s^t, a^t)$ . The learner updates the parameter  $\theta_l$  via policy gradient by computing the gradient  $\nabla_{\theta_l} f(\theta_l, \theta_e^*(\theta_l))$ , as the policy itself is parameterized by  $\theta_l$ .

**Overall learning process.** The entire learning procedure exhibits a hierarchical structure. Given the current learner’s reward parameter  $\theta_l$ , the lower level solves the ML-IRL problem (1) to determine the expert’s reward function  $r_{\theta_e^*(\theta_l)}$ . The upper level optimizes the expert-dependent reward function  $r_{\theta_l}$  to solve the RL problem (2). This hierarchical optimization framework is illustrated in Figure 1 and can be formulated as a bi-level optimization problem:

$$\arg \min_{\theta_l \in \Theta} f(\theta_l, \theta_e^*(\theta_l)), \quad \text{s.t.} \quad \theta_e^*(\theta_l) = \arg \min_{\theta_e \in \Theta} L(\theta_l, \theta_e). \quad (3)$$

Once  $r_{\theta_e^*(\theta_l^*)}$  and  $r_{\theta_l^*}$  are determined, the joint optimal policy  $\pi_{\theta_l^*, \theta_e^*(\theta_l^*)}$  is readily computed.

## 5 ALGORITHM

In this section, we develop a double-loop algorithm named GSIIRL (Algorithm 1) to solve the bi-level optimization problem described in (3). In each  $k$ -th outer loop iteration, the learner partially solves the lower-level optimization problem within an inner loop and subsequently employs this intermediate solution to tackle the upper-level optimization problem through an outer loop. The specifics of the inner and outer loops are detailed in Sections 5.1 and 5.2, respectively.

### 5.1 INNER LOOP

At each iteration  $t$  of the inner loop, the learner updates the parameter  $\theta_e(t)$  via projected SGD. Specifically,  $\theta_e(t)$  is updated by moving in the opposite direction of the partial gradient  $\nabla_{\theta_e} L(\theta_l(k), \theta_e(t))$  and projecting the resulting parameter back onto the feasible set  $\Theta$ . The analytical expression of the gradient  $\nabla_{\theta_e} L(\theta_l, \theta_e)$  is provided in the following lemma.

**Lemma 1.** *The gradient  $\nabla_{\theta_e} L(\theta_l, \theta_e) = \mu_e(\pi_{\theta_l, \theta_e}) - \hat{\mu}_e(D_{\theta_l, r_e}) + \lambda \theta_e$ , where  $\mu_e(\pi_{\theta_l, \theta_e}) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e} r_{\theta_e}(s^t, a^t)]$ ,  $\hat{\mu}_e(D_{\theta_l, r_e}) \triangleq \frac{1}{d} \sum_{i=1}^d \sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e} r_{\theta_e}(s^{it}, a^{it})$ , and  $(s^{it}, a^{it}) \in \zeta_i \in D_{\theta_l, r_e}$*

The inner loop uses  $t_k$ -step projected gradient descent  $\theta_e(t+1) = \Pi_{\Theta}(\theta_e(t) - \beta_t \nabla_{\theta_e} L(\theta_l(k), \theta_e(t)))$  to obtain an expert reward estimate  $r_{\theta_e(t_k-1)}$ .

### 5.2 OUTER LOOP

At  $k$ -th iteration of the outer loop, the learner updates the parameter  $\theta_l(k)$  via projected SGD. Similar to the inner loop, this update requires computing the hypergradient  $\nabla f(\theta_l(k), \theta_e(k))$ . The analytical form of  $\nabla f(\theta_l(k), \theta_e(k))$ , which is widely used in bi-level optimization problems, is presented in the following lemma.

**Lemma 2.** *The analytical expression of the hypergradient  $\nabla f(\theta_l, \theta_e)$  used for updating  $\theta_l$  is  $\nabla_{\theta_l} f(\theta_l, \theta_e) - \nabla_{\theta_l \theta_e}^2 L(\theta_l, \theta_e) [\nabla_{\theta_e}^2 L(\theta_l, \theta_e)]^{-1} \nabla_{\theta_e} f(\theta_l, \theta_e)$  (Xu & Zhu, 2023; Ghadimi & Wang, 2018; Colson et al., 2007).*

**Algorithm 1** General Scenario Interactive Inverse Reinforcement Learning (GSIIRL)

---

```

270 Initializes  $\theta_l(0) \in \Theta, \theta_e(0) \in \Theta$ , step size sequence  $\{\alpha_k\}, \{\beta_t\}$ , regularization parameter  $\lambda$  and
271 integer sequence  $\{t_k\}$ 
272
273 for  $k = 0, 1, \dots, K - 1$  do
274    $\theta_e(0) = \theta_e(k)$ 
275   Samples the trajectory set  $D_{\theta_l(k), r_e}$ 
276   for  $t = 0, \dots, t_k - 1$  do
277     Calculates  $\nabla_{\theta_e} L(\theta_l(k), \theta_e(t))$  following Lemma 1
278      $\theta_e(t + 1) = \Pi_{\Theta}(\theta_e(t) - \beta_t \nabla_{\theta_e} L(\theta_l(k), \theta_e(t)))$ 
279   end for
280    $\theta_e(k) = \theta_e(t_k - 1)$ 
281   Initializes the random vector  $\Delta(k)$  and the positive scalar  $p(k)$ 
282   Gets  $\hat{\nabla}_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k))$ ,  $\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))$ ,  $\hat{\nabla}_{\theta_l} f(\theta_l, \theta_e)$  and  $\hat{\nabla}_{\theta_e} f(\theta_l, \theta_e)$  through
283   SPSA approximation following equation (4)
284   Calculates  $[\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k))$  following equation (5)
285   Calculates  $\hat{\nabla} f(\theta_l(k), \theta_e(k))$  following equation (6)
286    $\theta_l(k + 1) = \Pi_{\Theta}(\theta_l(k) - \alpha_k \hat{\nabla} f(\theta_l(k), \theta_e(k)))$ 
287 end for

```

---

Computing  $\nabla f(\theta_l(k), \theta_e(k))$  involves the partial gradients  $\nabla_{\theta_l} f(\theta_l, \theta_e)$ ,  $\nabla_{\theta_e} f(\theta_l, \theta_e)$ , the Jacobian  $\nabla_{\theta_l \theta_e}^2 L(\theta_l, \theta_e)$  and the inverse Hessian  $[\nabla_{\theta_e}^2 L(\theta_l, \theta_e)]^{-1}$ . However, directly calculating these quantities has a large computational complexity of  $\mathcal{O}(T^2)$ . To address this issue, we adopt the SPSA method, which reduces computational complexity to  $\mathcal{O}(T)$ . A detailed discussion on the computational complexity reduction is provided later in Theorem 1. Using SPSA, the learner obtains the estimated hypergradient  $\hat{\nabla} f(\theta_l(k), \theta_e(k))$ , which serves as the gradient approximation in the projected SGD update. To calculate  $\hat{\nabla} f(\theta_l(k), \theta_e(k))$ , the learner estimates  $\hat{\nabla}_{\theta_l} f(\theta_l(k), \theta_e(k))$ ,  $\hat{\nabla}_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k))$ ,  $\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))$ , and  $\hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k))$ . Let us take  $\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))$  as an example to illustrate SPSA. According to the equation (2.2) in (Spall, 1992), it is approximated as follows:

$$\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k)) = \left[ \frac{\Delta_{d_k}}{2p\Delta_1(k)} \quad \dots \quad \frac{\Delta_{d_k}}{2p\Delta_m(k)} \right]^T, \quad (4)$$

where  $\Delta_{d_k} = \nabla_{\theta_e} L(\theta_l(k), \theta_e(k) + p(k)\Delta(k)) - \nabla_{\theta_e} L(\theta_l(k), \theta_e(k) - p(k)\Delta(k))$ , the perturbation  $\Delta(k) \in \mathbb{R}^m$  is a vector of  $m$  mutually independent zero-mean random variables and each element of  $\Delta(k)$  satisfies  $|\Delta_i(k)| \leq \alpha_0$ ,  $E|\Delta_i^{-1}(k)| \leq \alpha_1$ ,  $i = 1, \dots, m$  with  $\alpha_0, \alpha_1$  as positive constants. The parameter  $p(k)$  is a positive scalar.

Referring to equation (4), the computation of  $\hat{\nabla}_{\theta_l} f(\theta_l(k), \theta_e(k))$  requires  $f(\theta_l(k), \theta_e(k))$  with weight perturbations on  $\theta_l(k)$ . Analogously, computing  $\hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k))$ ,  $\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))$ , and  $\hat{\nabla}_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k))$  requires  $f(\theta_l(k), \theta_e(k))$ ,  $\nabla_{\theta_e} L(\theta_l(k), \theta_e(k))$ , and  $\nabla_{\theta_l} L(\theta_l(k), \theta_e(k))$  with weight perturbations on  $\theta_e(k)$ , respectively. Notice that the analytical expression of  $f(\theta_l(k), \theta_e(k))$  and  $\nabla_{\theta_e} L(\theta_l(k), \theta_e(k))$  are shown in optimization problem (2) and Lemma 1 respectively. Similarly to Lemma 1, the gradient  $\nabla_{\theta_l} L(\theta_l(k), \theta_e(k)) \triangleq \mu_l(\pi_{\theta_l(k), \theta_e(k)}) - \hat{\mu}_l(D_{\theta_l(k), r_e})$  and the proof is given in the Appendix A.4. The expectation  $\mu_l(\pi_{\theta_l, \theta_e}) \triangleq E^{\pi_{\theta_l, \theta_e}}[\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_l} r_{\theta_l}(s^t, a^t)]$  is the reward gradient expectation of the learner and the empirical expectation  $\hat{\mu}_l(D_{\theta_l, r_e}) \triangleq \frac{1}{d} \sum_{i=1}^d \sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_l} r_{\theta_l}(s^{it}, a^{it}), (s^{it}, a^{it}) \in \zeta_i \in D_{\theta_l, r_e}$  is the estimated reward gradient expectation of the learner.

To compute the hypergradient  $\hat{\nabla} f(\theta_l(k), \theta_e(k))$ , we must ensure that the estimated Hessian  $\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))$  is invertible. Since the negative log-likelihood function is convex and the  $L_2$  regularization term is strongly convex, the overall function  $L(\theta_l(k), \theta_e(k))$  is strongly convex. Hence its Hessian is positive definite. By choosing suitable  $\Delta(k)$  and  $p(k)$ , we ensure that the estimated Hessian  $\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))$  is positive definite. Since directly inverting a matrix is com-

putationally expensive, we use the conjugate gradient method to compute the product

$$\begin{aligned} & [\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k)) \\ &= \min_u \frac{1}{2} u^T \hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k)) u - u^T \hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k)). \end{aligned} \quad (5)$$

Finally, using  $\hat{\nabla}_{\theta_l} f(\theta_l(k), \theta_e(k))$ ,  $\hat{\nabla}_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k))$ , and the result from conjugate gradient, we compute the estimated hypergradient

$$\begin{aligned} & \hat{\nabla} f(\theta_l(k), \theta_e(k)) \\ &= \hat{\nabla}_{\theta_l} f(\theta_l(k), \theta_e(k)) - \hat{\nabla}_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k)) [\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k)). \end{aligned} \quad (6)$$

After  $K$  iterations, the outer loop terminates, yielding  $r_{\theta_l(K)}$ ,  $r_{\theta_e(K)}$  and  $\pi_{\theta_l(K), \theta_e(K)}$ .

## 6 ANALYTICAL RESULT

In this section, we present our analytical results regarding computational complexity reduction and convergence rate. To facilitate our analysis, we impose the following assumption on the estimated reward functions  $r_{\theta_l}$  and  $r_{\theta_e}$ :

**Assumption 1.** *The estimated expert-dependent reward function  $r_{\theta_l}$  and the estimated expert reward function  $r_{\theta_e}$  are fourth differentiable, i.e.,  $C^4$ .*

Since  $\theta_l$  and  $\theta_e$  lie in compact sets, Assumption 1 implies that the derivatives of the reward functions  $r_{\theta_l}$  and  $r_{\theta_e}$  are bounded. Such boundedness of higher-order derivatives is a standard assumption in the literature and has been widely adopted in bi-level optimization (Jin et al., 2020; Zeng et al., 2022; Liu & Zhu, 2023a), RL (Wang et al., 2019; Zhang et al., 2020), and IRL (Liu & Zhu, 2023b).

### 6.1 COMPUTATIONAL COMPLEXITY

The computation complexities of computing  $\nabla f(\theta_l, \theta_e)$  and  $\hat{\nabla} f(\theta_l, \theta_e)$  are shown in Theorem 1.

**Theorem 1.** *Consider  $T$  as the decision factor, the computational complexity of computing  $\nabla f(\theta_l, \theta_e)$  is  $\mathcal{O}(T^2)$  and that of computing  $\hat{\nabla} f(\theta_l, \theta_e)$  is  $\mathcal{O}(T)$ .*

Applying SPSA to approximate  $\nabla f(\theta_l, \theta_e)$  reduces the per-iteration computational complexity of the upper-level problem from  $\mathcal{O}(T^2)$  to  $\mathcal{O}(T)$ . The proof of Theorem 1 is in the Appendix A.8. Compared to finite-difference methods, SPSA requires fewer policies for the approximation. In SPSA, a single random perturbation of  $\theta_l$  or  $\theta_e$  yields two policies (one for the positive perturbation and one for the negative perturbation). In contrast, the finite difference approach requires two policies for each parameter dimension. Each policy is generated via MARL, so the computational cost of producing these policies is non-negligible. As SPSA requires fewer policies, its overall computational cost of SPSA is lower than that of finite difference methods. Finite-difference methods, in turn, have a lower computational cost than explicitly computing gradients.

### 6.2 CONVERGENCE RATE

The convergence of our algorithm is shown in Theorem 2.

**Theorem 2.** *Suppose Assumption 1 holds, by choosing  $p(k) = \frac{1}{k}$ ,  $\alpha_k = \frac{1}{L_f \sqrt{K}}$ ,  $t_k = \lceil \frac{\sqrt{k+1}}{2} \rceil$ , the following convergence guarantee holds:  $\frac{1}{K} \sum_{k=0}^{K-1} E[\|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2] \leq \mathcal{O}(\frac{1}{\sqrt{K}})$ .*

Theorem 2 indicates that the expected hypergradient decays at a rate of  $\mathcal{O}(\frac{1}{\sqrt{K}})$ , matching the convergence rate of standard bilevel optimization (Ghadimi & Wang, 2018). This implies that the bias introduced by SPSA does not slow convergence compared to standard bilevel optimization. Corollary 2.1 shows that the linear reward function  $r_{\theta_e}$  is a sufficient condition for the convergence of the policy  $\pi_{\theta_e}$ . Convergence of the cumulative reward difference has been widely used to infer convergence of the learned expert policy (Rhinehart & Kitani, 2017; Renard et al., 2024).

**Corollary 2.1.** *If the reward function  $r_{\theta_e}$  is linear, the cumulative reward difference between the learned expert policy  $\pi_{\theta_e}$  and  $\pi_e$  decreases at a rate  $\mathcal{O}(\frac{1}{\sqrt{K}})$ .*

## 7 EXPERIMENT

In our experiments, we evaluate GSIIRL on four widely used environments. The Multi-Agent Particle Environment (MPE) (Lowe et al., 2017; Terry et al., 2021; Yu et al., 2022) involves particle agents that can move, communicate, observe one another, push each other, and interact with fixed landmarks. The StarCraft Multi-Agent Challenge (SMAC) (Samvelyan et al., 2019; Yu et al., 2022; Rashid et al., 2020) is designed based on the real-time strategy game StarCraft II. The Human-Robot Interaction environment (El-Shamouty et al., 2020; Fan et al., 2020; Zhu & Zhang, 2021) models a navigation scenario in which a robot interacts with a human; the corresponding results are reported in Appendix Section A.9.3. The security environment (Zhang et al., 2019; 2021) involves automated defense systems; the results are presented in Appendix Section A.9.4.

We compare GSIIRL with four baseline algorithms: the **MARL** algorithm (Lowe et al., 2017), the **MA-IRL** algorithm (Ziebart et al., 2010), the cooperative interactive IRL (**CIRL**) algorithm (Büning et al., 2022), and the maximum likelihood IRL (**ML-IRL**) algorithm (Zeng et al., 2022). Since none of these baseline algorithms can directly address the interactive IRL problem, we adapt the experimental setup to make them applicable. Specifically, for MARL, we assume that both the learner and the expert know their reward functions. For MA-IRL, we provide demonstrations sampled according to ground-truth reward functions of both the learner and the expert. For CIRL, we assume that the expert-dependent reward function of the learner is identical to the expert’s reward function. For ML-IRL, we assume that the learner has access to demonstrations generated by a policy through MARL with the learner’s initial and the ground-truth expert reward functions.

The MARL baseline serves to illustrate the best achievable performance among all related algorithms. The MA-IRL baseline evaluates whether our algorithm can reach similar performance with less information. The CIRL baseline demonstrates that cooperative interactive IRL algorithms are only suitable for fully cooperative cases of interactive IRL problems. The ML-IRL baseline shows that conventional IRL algorithms are not well-suited for interactive IRL problems.

### 7.1 MULTI-AGENT PARTICLE ENVIRONMENTS

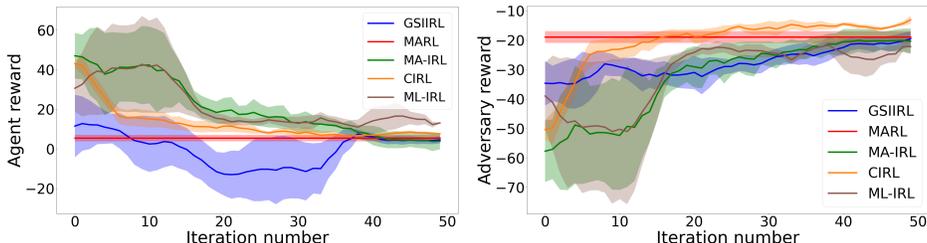


Figure 2: MPE simulation results. **Left:** Agent’s reward. **Right:** Adversary’s reward. The horizontal lines show the cumulative rewards from a MARL method that has access to the ground-truth reward functions after convergence. In contrast, other methods learn the reward functions from demonstrations and interactions, respectively. In each iteration, we compute the policies of both the adversary and the agents via MARL using the currently learned reward functions. We then evaluate these policies using the ground-truth reward functions.

We use the physical deception environment in the MPE for our simulation. This environment includes one adversary, two cooperating agents, and two landmarks. One of the landmarks is the target landmark. The cooperating agents know which landmark is the target, whereas the adversary does not. The cooperating agents aim to have one of them approach the target as closely as possible while preventing the adversary from reaching it. Conversely, the adversary aims to identify and reach the target landmark. We treat the two cooperating agents as a single learner and the adversary as the expert. Both the learner and the expert have continuous state and action spaces. The learner’s state space is 10-dimensional and the expert’s is 8-dimensional, and each has a 5-dimensional action space. Other details are provided in Appendix Section A.9.1.

Because the two agents cooperate and share a common reward, we plot a single cumulative reward curve for both. Figure 2 presents results for two randomly initialized reward functions updated over 50 iterations. Initially, both the adversary and the agents act randomly, and the cumulative

rewards differ significantly from those of the MARL baseline. As more data are collected, the agents gradually learn their own goal and the adversary’s goal, leading the cumulative rewards of GSIIRL and MA-IRL to converge to that of MARL. In contrast, CIRL and ML-IRL fail to do so.

## 7.2 STARCRAFT MULTI-AGENT CHALLENGE

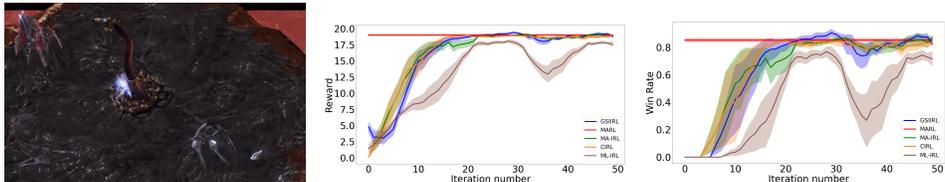


Figure 3: SMAC simulation results. **Left:** A screenshot of the game. **Middle:** Cumulative reward. **Right:** Win rate. The cumulative rewards are computed in the same manner as described in Figure 2. Similarly, the win rates are computed based on the policies corresponding to the learned reward functions. The agents are considered victorious when they defeat the enemy.

We use the scenario `2s_vs_1sc` from the SMAC for our experiment. In this scenario, two agents cooperate to defeat a single, more powerful enemy unit (which is controlled by a built-in AI). Because the environment does not expose the enemy’s state or action information, the enemy’s reward function cannot be learned directly. We therefore designate one agent as the expert (with access to the ground-truth reward function and the optimal cooperative strategy to defeat the enemy) and the other agent as the learner. The learner’s goal is to infer the expert’s strategy and learn to cooperate with the expert accordingly. Since this setting is fully cooperative, the agents share a common reward; accordingly, we ignore any expert-independent component of the learner’s reward and consider only the expert-dependent component. Both agents have a continuous 17-dimensional state space and a discrete action space of size 7. Additional implementation details are provided in Appendix A.9.2.

Because the two agents share the same reward, we report a single cumulative reward curve. The middle panel of Figure 3 shows that the cumulative rewards of GSIIRL and MA-IRL converge to that achieved by MARL, mirroring the trend observed in Figure 2. Since this scenario is fully cooperative, the cumulative reward of CIRL also converges to that of the MARL, whereas ML-IRL does not. We further evaluate performance using the agents’ win rate. At each iteration, we obtain the agents’ policies via MARL with the current learned reward functions. Each learned policy is then evaluated over 100 independent trials, and the win rate (the fraction of trials in which the agents defeat the enemy) is recorded. The right panel of Figure 3 shows that the win rates of GSIIRL, MA-IRL, and CIRL converge to that of the MARL baseline, while ML-IRL fails to do so.

## 7.3 RESULTS ANALYSIS

From Figures 2 and 3, we conclude that CIRL is effective only in fully cooperative cases, and ML-IRL is not suitable for interactive IRL problems. We therefore focus on comparing the cumulative rewards achieved by MA-IRL, GSIIRL, and MARL across the four experiments. Figure 2 and Figure 3 show that the cumulative rewards of MA-IRL and GSIIRL converge to those of MARL. The final cumulative rewards of GSIIRL are comparable to those of MA-IRL and MARL, differing by at most 5%, and a table with numerical cumulative reward comparison is shown in Appendix A.9.5. Notably, GSIIRL requires less information than MA-IRL: it relies only on the expert’s interactions with the environment and the corresponding values received during the interactions. In contrast, MA-IRL requires demonstrations generated using the ground-truth reward functions of both the learner and the expert, whereas MARL assumes full knowledge of the ground-truth reward functions.

## 8 CONCLUSION

We develop a general interactive IRL framework that enables a learner to infer an expert’s reward function and learn an appropriate interaction policy through interactions with the expert. This framework releases the strong assumption on the learner’s and the expert’s reward functions. We propose the GSIIRL algorithm and provide a theoretical analysis of its convergence rate. Experiments in both continuous and discrete environments demonstrate the effectiveness of GSIIRL.

## 486 9 ETHICS STATEMENT

487

488 This work does not involve human subjects or sensitive personal data. All experiments are conducted  
489 on publicly available simulation benchmarks.

490

491

## 492 10 REPRODUCIBILITY STATEMENT

493

494 The details of the experimental setup are provided in Appendix A.9, and the source code is included  
495 in the supplementary material to facilitate reproducibility.

496

## 497 REFERENCES

498

499 Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In  
500 *Proceedings of the twenty-first international conference on Machine learning*, 2004.

501

502 Zoe Ashwood, Aditi Jha, and Jonathan W Pillow. Dynamic inverse reinforcement learning for  
503 characterizing animal behavior. *Advances in neural information processing systems*, 35:29663–  
504 29676, 2022.

505

506 Thomas Kleine Büning, Anne-Marie George, and Christos Dimitrakakis. Interactive inverse rein-  
507 forcement learning for cooperative games. In *International Conference on Machine Learning*, pp.  
2393–2413. PMLR, 2022.

508

509 Jaedeug Choi and Kee-Eung Kim. Nonparametric bayesian inverse reinforcement learning for mul-  
510 tiple reward functions. *Advances in neural information processing systems*, 25, 2012.

511

512 Benoît Colson, Patrice Marcotte, and Gilles Savard. An overview of bilevel optimization. *Annals of*  
*operations research*, 153(1):235–256, 2007.

513

514 Mohamed El-Shamouty, Xinyang Wu, Shanqi Yang, Marcel Albus, and Marco F Huber. Towards  
515 safe human-robot collaboration using deep reinforcement learning. In *2020 IEEE international*  
516 *conference on robotics and automation (ICRA)*, pp. 4899–4905. IEEE, 2020.

517

518 Mahmoud Elnaggar and Nicola Bezzo. An irl approach for cyber-physical attack intention prediction  
and recovery. In *2018 Annual American Control Conference (ACC)*, pp. 222–227. IEEE, 2018.

519

520 Tingxiang Fan, Pinxin Long, Wenxi Liu, and Jia Pan. Distributed multi-robot collision avoidance  
521 via deep reinforcement learning for navigation in complex scenarios. *The International Journal*  
522 *of Robotics Research*, 39(7):856–892, 2020.

523

524 Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint*  
*arXiv:1802.02246*, 2018.

525

526 Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with  
527 deep energy-based policies. In *International conference on machine learning*, pp. 1352–1361.  
528 PMLR, 2017.

529

530 Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. Cooperative inverse  
531 reinforcement learning. *Advances in neural information processing systems*, 29, 2016.

532

533 Tsubasa Hirakawa, Takayoshi Yamashita, Toru Tamaki, Hironobu Fujiyoshi, Yuta Umezu, Ichiro  
534 Takeuchi, Sakiko Matsumoto, and Ken Yoda. Can ai predict animal movements? filling gaps in  
animal trajectories using inverse reinforcement learning. *Ecosphere*, 9(10), 2018.

535

536 Chi Jin, Praneeth Netrapalli, and Michael Jordan. What is local optimality in nonconvex-nonconcave  
537 minimax optimization? In *International conference on machine learning*, pp. 4880–4889. PMLR,  
538 2020.

539

Parameswaran Kamalaruban, Rati Devidze, Volkan Cevher, and Adish Singla. Interactive teaching  
algorithms for inverse reinforcement learning. *arXiv preprint arXiv:1905.11867*, 2019.

- 540 Charles D Kolstad and Leon S Lasdon. Derivative evaluation and computational experience with  
541 large bilevel mathematical programs. *Journal of optimization theory and applications*, 65:485–  
542 499, 1990.
- 543
- 544 Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong  
545 Yang. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv preprint*  
546 *arXiv:2109.11251*, 2021.
- 547
- 548 Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with  
549 differentiable convex optimization. In *Proceedings of the IEEE/CVF conference on computer*  
550 *vision and pattern recognition*, pp. 10657–10665, 2019.
- 551
- 552 Xiaomin Lin, Stephen C Adams, and Peter A Beling. Multi-agent inverse reinforcement learning  
553 for certain general-sum stochastic games. *Journal of Artificial Intelligence Research*, 66:473–502,  
554 2019.
- 555
- 556 Shicheng Liu and Minghui Zhu. Distributed inverse constrained reinforcement learning for multi-  
557 agent systems. *Advances in Neural Information Processing Systems*, 35:33444–33456, 2022.
- 558
- 559 Shicheng Liu and Minghui Zhu. Learning multi-agent behaviors from distributed and streaming  
560 demonstrations. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a.
- 561
- 562 Shicheng Liu and Minghui Zhu. Meta inverse constrained reinforcement learning: Convergence  
563 guarantee and generalization analysis. In *The Twelfth International Conference on Learning Rep-*  
564 *resentations*, 2023b.
- 565
- 566 Shicheng Liu and Minghui Zhu. In-trajectory inverse reinforcement learning: Learn incrementally  
567 before an ongoing trajectory terminates. *Advances in Neural Information Processing Systems*, 37:  
568 117164–117209, 2024.
- 569
- 570 Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-  
571 agent actor-critic for mixed cooperative-competitive environments. *Advances in neural informa-*  
572 *tion processing systems*, 30, 2017.
- 573
- 574 Farzan Memarian, Wonjoon Goo, Rudolf Lioutikov, Scott Niekum, and Ufuk Topcu. Self-supervised  
575 online reward shaping in sparse-reward environments. In *2021 IEEE/RSJ International Confer-*  
576 *ence on Intelligent Robots and Systems (IROS)*, pp. 2369–2375. IEEE, 2021.
- 577
- 578 Andrew Y Ng and Stuart Russell. Algorithms for inverse reinforcement learning. In *International*  
579 *Conference on Machine Learning*, 2000.
- 580
- 581 Billy Okal and Kai O Arras. Learning socially normative robot navigation behaviors with bayesian  
582 inverse reinforcement learning. In *2016 IEEE international conference on robotics and automa-*  
583 *tion (ICRA)*, pp. 2889–2895. IEEE, 2016.
- 584
- 585 Malayandi Palaniappan, Dhruv Malik, Dylan Hadfield-Menell, Anca Dragan, and Stuart Russell.  
586 Efficient cooperative inverse reinforcement learning. In *Proc. ICML Workshop on Reliable Ma-*  
587 *chine Learning in the Wild*, 2017.
- 588
- 589 Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In *International con-*  
590 *ference on machine learning*, pp. 737–746. PMLR, 2016.
- 591
- 592 Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. In *IJCAI*, vol-  
593 *ume 7*, pp. 2586–2591, 2007.
- 594
- 595 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster,  
596 and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement  
597 learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.
- 598
- 599 Titouan Renard, Andreas Schlaginhausen, Tingting Ni, and Maryam Kamgarpour. Convergence  
600 of a model-free entropy-regularized inverse reinforcement learning algorithm. *arXiv preprint*  
601 *arXiv:2403.16829*, 2024.

- 594 Nicholas Rhinehart and Kris M Kitani. First-person activity forecasting with online inverse rein-  
595 forcement learning. In *Proceedings of the IEEE International Conference on Computer Vision*,  
596 pp. 3696–3705, 2017.
- 597 Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas  
598 Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson.  
599 The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- 600 James C Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient  
601 approximation. *IEEE transactions on automatic control*, 37(3):332–341, 1992.
- 602 John C Strikwerda. *Finite difference schemes and partial differential equations*. SIAM, 2004.
- 603 Jordan Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan,  
604 Luis S Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. Petting-  
605 zoo: Gym for multi-agent reinforcement learning. *Advances in Neural Information Processing*  
606 *Systems*, 34:15032–15043, 2021.
- 607 Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global  
608 optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2019.
- 609 Xingyu Wang and Diego Klabjan. Competitive multi-agent inverse reinforcement learning with  
610 sub-optimal demonstrations. In *International Conference on Machine Learning*, pp. 5143–5151.  
611 PMLR, 2018.
- 612 Siyuan Xu and Minghui Zhu. Meta value learning for fast policy-centric optimal motion planning.  
613 In *Robotics science and systems*, 2022.
- 614 Siyuan Xu and Minghui Zhu. Efficient gradient approximation method for constrained bilevel op-  
615 timization. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(10):12509–12517,  
616 Jun. 2023.
- 617 Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The  
618 surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information*  
619 *Processing Systems*, 35:24611–24624, 2022.
- 620 Lantao Yu, Jiaming Song, and Stefano Ermon. Multi-agent adversarial inverse reinforcement learn-  
621 ing. In *International Conference on Machine Learning*, pp. 7194–7201. PMLR, 2019.
- 622 Siliang Zeng, Chenliang Li, Alfredo Garcia, and Mingyi Hong. Maximum-likelihood inverse rein-  
623 forcement learning with finite-time guarantees. *Advances in Neural Information Processing*  
624 *Systems*, 35:10122–10135, 2022.
- 625 Siliang Zeng, Chenliang Li, Alfredo Garcia, and Mingyi Hong. When demonstrations meet genera-  
626 tive world models: A maximum likelihood framework for offline inverse reinforcement learning.  
627 *Advances in Neural Information Processing Systems*, 36:65531–65565, 2023.
- 628 Dan Zhang, Gang Feng, Yang Shi, and Dipti Srinivasan. Physical safety and cyber security analysis  
629 of multi-agent systems: A survey of recent advances. *IEEE/CAA Journal of Automatica Sinica*, 8  
630 (2):319–333, 2021.
- 631 Kaiqing Zhang, Alec Koppel, Hao Zhu, and Tamer Basar. Global convergence of policy gradient  
632 methods to (almost) locally optimal policies. *SIAM Journal on Control and Optimization*, 58(6):  
633 3586–3612, 2020.
- 634 Xiangyuan Zhang, Kaiqing Zhang, Erik Miehling, and Tamer Basar. Non-cooperative inverse rein-  
635 forcement learning. *Advances in neural information processing systems*, 32, 2019.
- 636 Kai Zhu and Tao Zhang. Deep reinforcement learning based mobile robot navigation: A review.  
637 *Tsinghua Science and Technology*, 26(5):674–691, 2021.
- 638 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse  
639 reinforcement learning. *Association for the Advancement of Artificial Intelligence*, 8:1433–1438,  
640 2008.
- 641 Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling interaction via the principle of  
642 maximum causal entropy. *International Conference on Machine Learning*, 2010.
- 643  
644  
645  
646  
647

## A APPENDIX

### A.1 NOTION AND NOTATIONS

Define  $f(\theta_l, \theta_e) \triangleq J_{ed}(\pi_{\theta_l, \theta_e}) + J_{ei}(\pi_{\theta_l, \theta_e})$ , where  $J_{ed}(\pi_{\theta_l, \theta_e}) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t r_{\theta_l}(s^t, a^t)]$  is the cumulative expert-dependent reward of the learner and  $J_{ei}(\pi_{\theta_l, \theta_e}) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t r_{li}(s_l^t, a_l^t)]$  is the cumulative expert-independent reward of the learner. Define the reward gradient expectation of the expert from the state-action pair  $(s, a)$  as  $\mu_e(s, a) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e} r_e^{\theta_e}(s^t, a^t) | s^0 = s, a^0 = a]$ , the reward gradient expectation of the expert from the state  $s$  as  $\mu_e(s) \triangleq \int_{a_l \in a_l} \int_{a_e \in a_e} \mu_e(s, a) da_l da_e$ ,  $\mu_e(s, a) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e} r_e^{\theta_e}(s^t, a^t) | s^0 = s]$ . Analogously, define the reward gradient expectation of the learner as  $\mu_l(\pi_{\theta_l, \theta_e}) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_l} r_{\theta_l}(s^t, a^t)]$ , the reward gradient expectation of the learner from the state-action pair  $(s, a)$  as  $\mu_l(s) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_l} r_{\theta_l}(s^t, a^t) | s^0 = s]$ ,  $\mu_l(s, a) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_l} r_{\theta_l}(s^t, a^t) | s^0 = s, a^0 = a]$ , the reward gradient expectation of the learner from the state  $s$  as  $\mu_l(s) \triangleq \int_{a_l \in a_l} \int_{a_e \in a_e} \mu_l(s, a) da_l da_e$ . Define the cumulative expert-independent reward of the learner from the state-action pair  $(s, a)$  as  $J_{ei}(s, a) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t r_{li}(s_l^t, a_l^t) | s^0 = s, a^0 = a]$  and define the cumulative expert-independent reward of the learner from the state  $s$  as  $J_{ei}(s) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t r_{li}(s_l^t, a_l^t) | s^0 = s]$ . Define  $P_0(s)$  as the probability distribution of the initial state.

During the proofs, symbols  $(i)$  to  $(viii)$  are used to represent what theorems or methods are used to get the current step. The symbol  $(i)$  represents chain rule, the symbol  $(ii)$  represents the linearity of expectation, the symbol  $(iii)$  represents the close form of geometric series, the symbol  $(iv)$  represents the triangle inequality, the symbol  $(v)$  represents the Hölder's inequality, the symbol  $(vi)$  represents Taylor's theorem, the symbol  $(vii)$  means the usage of other equations in this paper, and the symbol  $(viii)$  means keeping expansion.

### A.2 FUNDAMENTAL RESULT FOR POLICY

This section lists the expressions for important gradients in continuous state-action space. Based on the idea of the soft Q learning (Haarnoja et al., 2017), we can get:

$$Q_{\theta_l, \theta_e}^{soft}(s, a) = r_{\theta_l}(s, a) + r_{li}(s, a) + r_e^{\theta_e}(s, a) + \gamma \int_{s' \in S} P(s' | s, a) V^{soft}(s') ds', \quad (7)$$

$$V_{\theta_l, \theta_e}^{soft}(s) = \ln \int_{a_l \in a_l} \int_{a_e \in a_e} \exp(Q^{soft}(s, a)) da_l da_e, \quad (8)$$

$$\pi_{\theta_l, \theta_e}(a_l, a_e | s) = \frac{\exp(Q^{soft}(s, a))}{\exp(V^{soft}(s))}. \quad (9)$$

As the lower-level optimization problem is a maximum likelihood problem, it is important to know  $\nabla_{\theta_e} \ln(\pi_{\theta_l, \theta_e})$  for the SGD of the lower-level optimization problem. The process for getting  $\nabla_{\theta_e} \ln(\pi_{\theta_l, \theta_e})$  is shown below.

From the expression of  $\pi_{\theta_l, \theta_e}$ , we can get  $\nabla_{\theta_e} \ln(\pi_{\theta_l, \theta_e}) = \nabla_{\theta_e} Q^{soft}(s, a) - \nabla_{\theta_e} V^{soft}(s)$ , therefore, we can calculate  $\nabla_{\theta_e} Q^{soft}(s, a)$  and  $\nabla_{\theta_e} V^{soft}(s)$  separately.

Based on the equation of  $V^{soft}(s)$ , the gradient  $\nabla_{\theta_e} V^{soft}(s)$  could be calculated as follows:

$$\begin{aligned}
& \nabla_{\theta_e} V^{soft}(s), \\
& \stackrel{(i)}{=} \frac{\int_{a_l \in a_l} \int_{a_e \in a_e} \nabla_{\theta_e} \exp(Q^{soft}(s, a)) da_l da_e}{\int_{a_l \in a_l} \int_{a_e \in a_e} \exp(Q^{soft}(s, a)) da_l da_e}, \\
& \stackrel{(vii)}{=} \frac{\int_{a_l \in a_l} \int_{a_e \in a_e} \exp(Q^{soft}(s, a)) \nabla_{\theta_e} Q^{soft}(s, a) da_l da_e}{\exp(V^{soft}(s))}, \\
& \stackrel{(vii)}{=} \int_{a_l \in a_l} \int_{a_e \in a_e} \pi_{\theta_l, \theta_e}(a_l, a_e | s) (\nabla_{\theta_e} r_e(s, a) \\
& + \gamma \int_{s' \in S} P(s' | s, a) \nabla_{\theta_e} V^{soft}(s') ds') da_l da_e, \\
& \stackrel{(viii)}{=} \int_{a_l \in a_l} \int_{a_e \in a_e} \pi_{\theta_l, \theta_e}(a_l, a_e | s) (\nabla_{\theta_e} r_e(s, a) + \gamma \int_{s' \in S} P(s' | s, a) \\
& (\int_{a'_l \in a_l} \int_{a'_e \in a_e} \pi_{\theta_l, \theta_e}(a'_l, a'_e | s') [\nabla_{\theta_e} r_e(s', a'_l, a'_e) \\
& + \gamma \int_{s'' \in S} P(s'' | s', a'_l, a'_e) \nabla_{\theta_e} V^{soft}(s'') ds''] da'_l da'_e ds') da_l da_e, \\
& = E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e} r_e(s^t, a^t) | s^0 = s], \\
& = \mu_e(s), \tag{10}
\end{aligned}$$

where the first (vii) uses the expression in equation (8) and (9), and the second (viii) uses the expression in equation (7).

Similarly, we can get  $\nabla_{\theta_e} Q^{soft}(s, a)$  from  $Q^{soft}(s, a)$ .

$$\begin{aligned}
& \nabla_{\theta_e} Q^{soft}(s, a), \\
& \stackrel{(vii)}{=} \nabla_{\theta_e} r_e(s, a) + \gamma \int_{s' \in S} P(s' | s, a) \nabla_{\theta_e} V^{soft}(s') ds' da_l da_e, \\
& \stackrel{(viii)}{=} E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e} r_e(s^t, a^t) | s^0 = s, a^0 = a], \\
& = \mu_e(s, a), \tag{11}
\end{aligned}$$

where (vii) uses the expression in equation (7).

By summing the results of equation (10) and (11), we can get  $\nabla_{\theta_e} \ln(\pi_{\theta_l, \theta_e})$  as follows:

$$\nabla_{\theta_e} \ln(\pi_{\theta_l, \theta_e}(a_l, a_e | s)) = \nabla_{\theta_e} Q^{soft}(s, a) - \nabla_{\theta_e} V^{soft}(s) = \mu_e(s, a) - \mu_e(s). \tag{12}$$

The way to get  $\nabla_{\theta_l} \ln(\pi_{\theta_l, \theta_e})$  is same as that of  $\nabla_{\theta_e} \ln(\pi_{\theta_l, \theta_e})$ , and the gradient  $\nabla_{\theta_l} \ln(\pi_{\theta_l, \theta_e}(a_l, a_e | s)) = \mu_l(s, a) - \mu_l(s)$ .

### A.3 THE PROOF OF LEMMA 1

In this section, we derived gradients that are necessary for our method.

Define  $P_D(s^t = s, a_l^t = a_l, a_e^t = a_e) \triangleq \begin{cases} 1 & s^t = s, a_l^t = a_l, a_e^t = a_e \\ 0 & \text{otherwise} \end{cases}$ , where  $(s, a) \in D_{\theta_l, r_e}$ , as the probability of  $(s, a)$  occurring at time  $t$  in the demonstration set  $D_{\theta_l, r_e}$ . With the fundamental

result of the policy, we can derive the  $\nabla_{\theta_e} L(\theta_l, \theta_e)$  and prove the Lemma 1 as follows:

$$\begin{aligned}
& \nabla_{\theta_e} L(\theta_l, \theta_e), \\
&= - \sum_{t=0}^{T-1} \gamma^t \int_{s \in S} \int_{a_l \in a_l} \int_{a_e \in a_e} P_D(s^t = s, a_l^t = a_l, a_e^t = a_e) \nabla_{\theta_e} \ln(\pi_{\theta_l, \theta_e}) da_e da_l ds + \lambda \theta_e, \\
&\stackrel{(vii)}{=} - \sum_{t=0}^{T-1} \gamma^t \int_{s \in S} \int_{a_l \in a_l} \int_{a_e \in a_e} P_D(s^t = s, a_l^t = a_l, a_e^t = a_e) (\nabla_{\theta_e} r_e(s, a) \\
&\quad + E^{\pi_{\theta_l, \theta_e}} [\sum_{t'=t+1}^{T-1} \gamma^{t'-t-1} \nabla_{\theta_e} r_e(s^t, a^t) | s^t = s, a_l^t = a_l, a_e^t = a_e] \\
&\quad - E^{\pi_{\theta_l, \theta_e}} [\sum_{t'=t}^{T-1} \gamma^{t'-t} \nabla_{\theta_e} r_e(s^t, a^t) | s^t = s]) da_e da_l ds + \lambda \theta_e, \\
&\stackrel{(ii)}{=} - \sum_{t=0}^{T-1} \gamma^t \int_{s \in S} \int_{a_l \in a_l} \int_{a_e \in a_e} P_D(s^t = s, a_l^t = a_l, a_e^t = a_e) (\nabla_{\theta_e} r_e(s, a) da_e da_l ds \\
&\quad - \int_{s \in S} \int_{a_l \in a_l} \int_{a_e \in a_e} P_D(s^0 = s, a_l^0 = a_l, a_e^0 = a_e) \\
&\quad E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e} r_e(s^t, a^t) | s^0 = s]) da_e da_l ds + \lambda \theta_e, \\
&= \mu_e(\pi_{\theta_l, \theta_e}) - \hat{\mu}_e(D_{\theta_l, r_e}) + \lambda \theta_e,
\end{aligned}$$

where (vii) uses the expression of equation (12).

#### A.4 OTHER NEEDED GRADIENTS

**Lemma 3.** Suppose Assumption 1 holds, the first-order gradients  $\nabla_{\theta_l} L(\theta_l, \theta_e) = \mu_l(\pi_{\theta_l, \theta_e}) - \hat{\mu}_l(D_{\theta_l, r_e})$ ,  $\nabla_{\theta_l} f(\theta_l, \theta_e) = E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t [(\mu_l(s^t, a^t) - \mu_l(s^t))(J_{ed}(s^t, a^t) + J_{ei}(s^t, a^t))^T]] + \mu_l(\pi_{\theta_l, \theta_e})$ ,  $\nabla_{\theta_e} f(\theta_l, \theta_e) = E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t [(\mu_e(s^t, a^t) - \mu_e(s^t))(J_{ed}(s^t, a^t) + J_{ei}(s^t, a^t))^T]]$ . The second-order gradients  $\nabla_{\theta_l \theta_e}^2 L(\theta_l, \theta_e) = E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t (\mu_e(s, a) - \mu_e(s)) \mu_l(s, a)^T]$ ,  $\nabla_{\theta_e \theta_e}^2 L(\theta_l, \theta_e) = E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t [(\mu_e(s^t, a^t) - \mu_e(s^t)) \mu_e(s^t, a^t)^T + \nabla_{\theta_e}^2 r_e^{\theta_e}(s^t, a^t)]] - \frac{1}{d} \sum_{i=0}^d \sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e} r_e^{\theta_e}(s^{it}, a^{it}) + \lambda$

*Proof.* Through the same process of calculating  $\nabla_{\theta_e} L(\theta_l, \theta_e)$ , we can get the gradient  $\nabla_{\theta_l} L(\theta_l, \theta_e)$  following the same process in Lemma 1, where  $\mu_l(\pi_{\theta_l, \theta_e}) \triangleq E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_l} r_{ld}(s^t, a^t)]$ .

The process for getting  $\nabla_{\theta_l} f(\theta_l, \theta_e)$  is shown below.

From the equation of  $f(\theta_l, \theta_e)$ , we can see that  $\nabla_{\theta_l} f(\theta_l, \theta_e) = \nabla_{\theta_l} J_{ed}(\pi_{\theta_l, \theta_e}) + \nabla_{\theta_l} J_{ei}(\pi_{\theta_l, \theta_e})$ , then  $\nabla_{\theta_l} J_{ed}(\pi_{\theta_l, \theta_e})$  and  $\nabla_{\theta_l} J_{ei}(\pi_{\theta_l, \theta_e})$  can be calculated separately.

The calculation of deriving  $\nabla_{\theta_l} J_{ed}(\pi_{\theta_l, \theta_e})$  is as follows:

$$\begin{aligned}
& \nabla_{\theta_l} J_{ed}(\pi_{\theta_l, \theta_e}), \\
& \stackrel{(i)}{=} \int_{s^0 \in S} P_0(s^0) \int_{a_l^0 \in a_l} \int_{a_e^0 \in a_e} [\nabla_{\theta_l} \pi(a_l^0, a_e^0 | s^0) J_{ed}(s^0, a_l^0, a_e^0)^T \\
& + \pi(a_l^0, a_e^0 | s^0) \nabla_{\theta_l} J_{ed}(s^0, a_l^0, a_e^0)^T] da_l^0 da_e^0 ds^0, \\
& = \int_{s^0 \in S} P_0(s^0) \int_{a_l^0 \in a_l} \int_{a_e^0 \in a_e} [\nabla_{\theta_l} \pi(a_l^0, a_e^0 | s^0) J_{ed}(s^0, a_l^0, a_e^0)^T + \pi(a_l^0, a_e^0 | s^0) (\nabla_{\theta_l} r_{\theta_l}(s^0, a_l^0, a_e^0) \\
& + \gamma \int_{s^1 \in S} P(s^1 | s^0, a_l^0, a_e^0) \nabla_{\theta_l} J_{ed}(s^1) ds^1] da_l^0 da_e^0 ds^0, \\
& \stackrel{(viii)}{=} \int_{s^0 \in S} P_0(s^0) \int_{a_l^0 \in a_l} \int_{a_e^0 \in a_e} \{ \nabla_{\theta_l} \pi(a_l^0, a_e^0 | s^0) J_{ed}(s^0, a_l^0, a_e^0)^T + \pi(a_l^0, a_e^0 | s^0) [\nabla_{\theta_l} r_{\theta_l}(s^0, a_l^0, a_e^0) \\
& + \gamma \int_{s^1 \in S} P(s^1 | s^0, a_l^0, a_e^0) \int_{a_l^1 \in a_l} \int_{a_e^1 \in a_e} \nabla_{\theta_l} \pi(a_l^1, a_e^1 | s^1) J_{ed}(s^1, a_l^1, a_e^1)^T \\
& + \pi(a_l^1, a_e^1 | s^1) (\nabla_{\theta_l} r_{\theta_l}(s^1, a_l^1, a_e^1) + \gamma \int_{s^2 \in S} P(s^2 | s^1, a_l^1, a_e^1) \nabla_{\theta_l} J_{ed}(s^2) ds^2] ds^1 \} da_l^0 da_e^0 ds^0, \\
& = E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t \left( \frac{\nabla_{\theta_l} \pi(a_l^t, a_e^t | s^t)}{\pi(a_l^t, a_e^t | s^t)} J_{ed}(s^t, a^t)^T + \nabla_{\theta_l} r_{\theta_l}(s^t, a^t) \right) \right], \\
& = E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t (\nabla_{\theta_l} \ln(\pi_{\theta_l, \theta_e}) J_{ed}(s^t, a^t)^T + \nabla_{\theta_l} r_{\theta_l}(s^t, a^t)) \right], \\
& \stackrel{(vii)}{=} E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t [(\mu_l(s^t, a^t) - \mu_l(s^t)) J_{ed}(s^t, a^t)^T] + \mu_l(\pi_{\theta_l, \theta_e}) \right],
\end{aligned}$$

where (vii) use the expression of the equation (12).

Similarly, we can get  $\nabla_{\theta_l} J_{ei}(\pi_{\theta_l, \theta_e})$  as follows:

$$\begin{aligned}
& \nabla_{\theta_l} J_{ei}(\pi_{\theta_l, \theta_e}), \\
& \stackrel{(i)}{=} \int_{s^0 \in S} P_0(s^0) \int_{a_l^0 \in a_l} \int_{a_e^0 \in a_e} [\nabla_{\theta_l} \pi(a_l^0, a_e^0 | s^0) J_{ei}(s^0, a_l^0, a_e^0) \\
& + \pi(a_l^0, a_e^0 | s^0) (\nabla_{\theta_l} J_{ei}(s^0, a_l^0, a_e^0))] da_l^0 da_e^0 ds^0, \\
& \stackrel{(viii)}{=} \int_{s^0 \in S} P_0(s^0) \int_{a_l^0 \in a_l} \int_{a_e^0 \in a_e} [\pi(a_l^0, a_e^0 | s^0) \nabla_{\theta_l} \ln(\pi(a_l^0, a_e^0 | s^0)) J_{ei}(s^0, a_l^0, a_e^0) \\
& + \pi(a_l^0, a_e^0 | s^0) \gamma \int_{s^1 \in S} P(s^1 | s^0, a_l^0, a_e^0) \nabla_{\theta_l} J_{ei}(s^1) ds^1] da_l^0 da_e^0 ds^0, \\
& \stackrel{(vii)}{=} E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t (\mu_l(s^t, a^t) - \mu_l(s^t)) J_{ei}(s^t, a^t)^T \right],
\end{aligned}$$

where (vii) use the expression of the equation (12).

By summing the result of  $\nabla_{\theta_l} J_{ed}(\pi_{\theta_l, \theta_e})$  and  $\nabla_{\theta_l} J_{ei}(\pi_{\theta_l, \theta_e})$ , the result of  $\nabla_{\theta_l} f(\theta_l, \theta_e)$  is as follows:

$$\begin{aligned}
& \nabla_{\theta_l} f(\theta_l, \theta_e), \\
& = E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t [(\mu_l(s^t, a^t) - \mu_l(s^t)) J_{ed}(s^t, a^t)^T] + \mu_l(\pi_{\theta_l, \theta_e}) \right] \\
& + E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t (\mu_l(s^t, a^t) - \mu_l(s^t)) J_{ei}(s^t, a^t)^T \right], \\
& \stackrel{(ii)}{=} E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t [(\mu_l(s^t, a^t) - \mu_l(s^t)) (J_{ed}(s^t, a^t) + J_{ei}(s^t, a^t))^T] + \mu_l(\pi_{\theta_l, \theta_e}) \right].
\end{aligned}$$

The  $\nabla_{\theta_e} f(\theta_l, \theta_e)$  is calculated in the same way as the  $\nabla_{\theta_l} f(\theta_l, \theta_e)$ .

The process for getting  $\nabla_{\theta_e}^2 L(\theta_l, \theta_e)$  is shown below.

As we proved in Lemma 1, the gradient  $\nabla_{\theta_e} L(\theta_l, \theta_e) = \mu_e(\pi_{\theta_l, \theta_e}) - \hat{\mu}_e(D_{\theta_l, r_e}) + \lambda\theta_e$ , as a result, we can take the derivative of each term separately.

The derivative of  $\mu_e(\pi_{\theta_l, \theta_e})$  w.r.t  $\theta_e$  is calculated as follows:

$$\begin{aligned}
& \nabla_{\theta_e} \mu_e(\pi_{\theta_l, \theta_e}), \\
& \stackrel{(i)}{=} \int_{s^0 \in S} P_0(s^0) \int_{a_l^0 \in a_l} \int_{a_e^0 \in a_e} [\nabla_{\theta_e} \pi(a_l^0, a_e^0 | s^0) \mu_e(s^0, a_l^0, a_e^0)^T \\
& + \pi(a_l^0, a_e^0 | s^0) \nabla_{\theta_e} \mu_e(s^0, a_l^0, a_e^0)^T] da_l^0 da_e^0 ds^0, \\
& = \int_{s^0 \in S} P_0(s^0) \int_{a_l^0 \in a_l} \int_{a_e^0 \in a_e} [\nabla_{\theta_l} \pi(a_l^0, a_e^0 | s^0) \mu_e(s^0, a_l^0, a_e^0)^T \\
& + \pi(a_l^0, a_e^0 | s^0) (\nabla_{\theta_e}^2 r_e^{\theta_e}(s^0, a_l^0, a_e^0) + \gamma \int_{s^1 \in S} P(s^1 | s^0, a_l^0, a_e^0) \nabla_{\theta_e} \mu_e(s^1)] da_l^0 da_e^0 ds^0, \\
& \stackrel{(viii)}{=} \int_{s^0 \in S} P_0(s^0) \int_{a_l^0 \in a_l} \int_{a_e^0 \in a_e} \{ \nabla_{\theta_l} \pi(a_l^0, a_e^0 | s^0) \mu_e(s^0, a_l^0, a_e^0)^T \\
& + \pi(a_l^0, a_e^0 | s^0) [\nabla_{\theta_e}^2 r_e^{\theta_e}(s^0, a_l^0, a_e^0) \\
& + \gamma \int_{s^1 \in S} P(s^1 | s^0, a_l^0, a_e^0) \int_{a_l^1 \in a_l} \int_{a_e^1 \in a_e} \nabla_{\theta_e} \pi(a_l^1, a_e^1 | s^1) \mu_e(s^1, a_l^1, a_e^1)^T \\
& + \pi(a_l^1, a_e^1 | s^1) (\nabla_{\theta_e}^2 r_e^{\theta_e}(s^1, a_l^1, a_e^1) + \gamma \int_{s^2 \in S} P(s^2 | s^1, a_l^1, a_e^1) \nabla_{\theta_e} \mu_e(s^2)] ds^1 \} da_l^0 da_e^0 ds^0, \\
& = E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t \left( \frac{\nabla_{\theta_e} \pi(a_l^t, a_e^t | s^t)}{\pi(a_l^t, a_e^t | s^t)} \mu_e(s^t, a^t)^T + \nabla_{\theta_e}^2 r_e^{\theta_e}(s^t, a^t) \right) \right], \\
& = E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t (\nabla_{\theta_e} \ln(\pi_{\theta_l, \theta_e}) \mu_e(s^t, a^t)^T + \nabla_{\theta_e}^2 r_e^{\theta_e}(s^t, a^t)) \right], \\
& \stackrel{(vii)}{=} E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t [(\mu_e(s^t, a^t) - \mu_e(s^t)) \mu_e(s^t, a^t)^T + \nabla_{\theta_e}^2 r_e^{\theta_e}(s^t, a^t)] \right],
\end{aligned}$$

where (vii) use the expression of the equation (12).

With the result of  $\nabla_{\theta_e} \mu_e(\pi_{\theta_l, \theta_e})$ , the derivative of  $\nabla_{\theta_e} L(\theta_l, \theta_e)$  w.r.t  $\theta_e$  is as follows:

$$\begin{aligned}
& \nabla_{\theta_e}^2 L(\theta_l, \theta_e) = \nabla_{\theta_e} (\nabla_{\theta_e} L(\theta_l, \theta_e)), \\
& = \nabla_{\theta_e} (\mu_e(\pi_{\theta_l, \theta_e}) - \hat{\mu}_e(D_{\theta_l, r_e}) + \lambda\theta_e), \\
& \stackrel{(ii)}{=} E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t [(\mu_e(s^t, a^t) - \mu_e(s^t)) \mu_e(s^t, a^t)^T + \nabla_{\theta_e}^2 r_e^{\theta_e}(s^t, a^t)] \right] \\
& - \frac{1}{d} \sum_{i=0}^d \sum_{t=0}^{T-1} \gamma^t \nabla_{\theta_e}^2 r_e^{\theta_e}(s^{it}, a^{it}) + \lambda.
\end{aligned}$$

Through the same process of calculating  $\nabla_{\theta_e} \mu_e(\pi_{\theta_l, \theta_e})$ , the result is as follows:

$$\nabla_{\theta_l \theta_e}^2 L(\theta_l, \theta_e) = E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t (\mu_e(s, a) - \mu_e(s)) \mu_l(s, a)^T \right]$$

□

#### A.5 PROPERTIES OF THE LOWER LEVEL OPTIMIZATION PROBLEM

**Lemma 4.** Suppose Assumption 1 holds, for any  $\theta_l \in \mathbb{R}^n$  and  $\theta_e \in \mathbb{R}^m$ ,  $L$  is continuously twice differentiable in  $(\theta_l, \theta_e)$ .

918 For any  $\bar{\theta}_1 \in \mathbb{R}^n$ ,  $\nabla_{\theta_e} L(\bar{\theta}_1, \theta_e)$  is Lipschitz continuous (w.r.t  $\theta_e$ ) with constant  $L_{L_{\theta_e}} > 0$ .  
 919 For any  $\bar{\theta}_1 \in \mathbb{R}^n$  and  $\bar{\theta}_2 \in \mathbb{R}^m$ , we have  $\|\nabla_{\theta_l \theta_e}^2 L(\bar{\theta}_1, \bar{\theta}_2)\| \leq C_{L_{\theta_l \theta_e}}$  for some constant  $C_{L_{\theta_l \theta_e}} > 0$ .  
 920 For any  $\bar{\theta}_1 \in \mathbb{R}^n$ ,  $\nabla_{\theta_l \theta_e}^2 L(\bar{\theta}_1, \theta_e)$  and  $\nabla_{\theta_e \theta_e}^2 L(\bar{\theta}_1, \theta_e)$  are Lipschitz continuous (w.r.t  $\theta_e$ ) with  
 921 constants  $L_{L_{\theta_l \theta_e}} > 0$  and  $L_{L_{\theta_e \theta_e}} > 0$ .  
 922 For any  $\bar{\theta}_2 \in \mathbb{R}^m$ ,  $\nabla_{\theta_l \theta_e}^2 L(\theta_l, \bar{\theta}_2)$  and  $\nabla_{\theta_e \theta_e}^2 L(\theta_l, \bar{\theta}_2)$  are Lipschitz continuous (w.r.t  $\theta_l$ ) with  
 923 constants  $\bar{L}_{L_{\theta_l \theta_e}} > 0$  and  $\bar{L}_{L_{\theta_e \theta_e}} > 0$

924  
 925  
 926  
 927 *Proof.* Suppose that  $h$  is a real-valued function defined and differentiable on an interval  $H \subset \mathbb{R}$ . If  
 928  $\|\nabla h\|$  is bounded on  $H$ , then  $h$  is a Lipschitz function on  $H$ . So we need to prove  $\nabla_{\theta_e \theta_e}^2 L(\theta_l, \theta_e)$  is  
 929 bounded. From Assumption 1, we can show that  $\exists R_{1g} > 0, \|\nabla r_{\theta_l}\| \leq R_{1g}$ .

$$930 \quad \|\mu_l(s)\| \leq E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t R_{1g} |s^0 = s| \right] \stackrel{(i)}{\leq} \frac{R_{1g}}{1 - \gamma},$$

931 where (i) uses the close form of geometric series.

932 As a result,  $\|\mu_l(s)\|$  is bounded, proved through the same way,  $\|\mu_e(s)\|, \|\mu_l(s, a)\|, \|\mu_e(s, a)\|$   
 933 are also bounded. Based on the Lemma 3, all elements of  $\nabla_{\theta_e \theta_e}^2 L(\theta_l, \theta_e)$  are finite, therefore,  
 934  $\|\nabla_{\theta_e \theta_e}^2 L(\theta_l, \theta_e)\|$  is bounded and the  $\nabla_{\theta_e} L(\bar{\theta}_1, \theta_e)$  is Lipschitz continuous.

935 In the same way of proving  $\|\nabla_{\theta_e \theta_e}^2 L(\theta_l, \theta_e)\|$  is bounded, we can show  $\|\nabla_{\theta_l \theta_e}^2 L(\theta_l, \theta_e)\|$  is bounded.

936  
 937 We need to prove  $\nabla_{\theta_e \theta_e \theta_e}^3 L(\theta_l, \theta_e)$  and  $\nabla_{\theta_l \theta_e \theta_e}^3 L(\theta_l, \theta_e)$  are bounded. The proof of  $\nabla_{\theta_l \theta_e \theta_e}^3 L(\theta_l, \theta_e)$   
 938 is bounded as follows:

$$939 \quad \begin{aligned} & \nabla_{\theta_l \theta_e \theta_e}^3 L(\theta_l, \theta_e), \\ &= \nabla_{\theta_e} \left( E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t (\mu_e(s, a) - \mu_e(s)) \mu_l(s, a)^T \right] \right), \\ & \stackrel{(i)}{=} E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t (\nabla_{\theta_e} \mu_e(s, a)) \mu_l(s, a)^T + \mu_e(s, a) (\nabla_{\theta_e} \mu_l(s, a)^T) \right. \\ & \quad \left. - (\nabla_{\theta_e} \mu_e(s)) \mu_l(s, a)^T - \mu_e(s) (\nabla_{\theta_e} \mu_l(s, a)^T) \right]. \end{aligned}$$

940 Each gradient inside the expectation could be derived through the process of deriving  $\nabla_{\theta_l} \mu_l(\pi_{\theta_l, \theta_e})$   
 941 in the proof of Lemma 3 and these gradients are all finite with the same way of proving  
 942  $\|\nabla_{\theta_e \theta_e}^2 L(\theta_l, \theta_e)\|$  is bounded. The third-order gradient  $\nabla_{\theta_l \theta_e \theta_e}^3 L(\theta_l, \theta_e)$  is bounded with the same  
 943 process. Therefore, the third-order gradients of  $L(\theta_l, \theta_e)$  are all bounded.

944 Through the same procedure, we can prove  $\nabla_{\theta_e \theta_e \theta_e}^3 L(\theta_l, \theta_e)$  and  $\nabla_{\theta_l \theta_e \theta_e}^3 L(\theta_l, \theta_e)$  are bounded.

945 At the same time, the existence of  $\nabla_{\theta_l \theta_e}^2 L(\theta_l, \theta_e), \nabla_{\theta_e \theta_e}^2 L(\theta_l, \theta_e)$  are shown. Analogously, the  
 946 existence of  $\nabla_{\theta_l \theta_l}^2 L(\theta_l, \theta_e), \nabla_{\theta_e \theta_l}^2 L(\theta_l, \theta_e)$  could be proved in the same way. The third-order  
 947 gradients of  $L(\theta_l, \theta_e)$  are bounded. Therefore,  $L$  is continuously twice differentiable in  $(\theta_l, \theta_e)$

948 □

## 949 A.6 PROPERTIES OF THE UPPER-LEVEL OPTIMIZATION PROBLEM

950 **Lemma 5.** Suppose Assumption 1 holds, for any  $\bar{\theta}_1 \in \mathbb{R}^n$ ,  $\nabla_{\theta_l} f(\bar{\theta}_1; \theta_e)$  and  $\nabla_{\theta_e} f(\bar{\theta}_1; \theta_e)$  are  
 951 Lipschitz continuous (w.r.t  $\theta_e$ ) with constants  $L_{f_{\theta_l}} > 0$  and  $L_{f_{\theta_e}} > 0$ .

952 For any  $\bar{\theta}_2 \in \mathbb{R}^m$ ,  $\nabla_{\theta_e} f(\theta_l; \bar{\theta}_2)$  is Lipschitz continuous (w.r.t  $\theta_l$ ) with constants  $\bar{L}_{f_{\theta_e}} > 0$

953 For any  $\bar{\theta}_1 \in \mathbb{R}^n$  and  $\bar{\theta}_2 \in \mathbb{R}^m$ , we have  $\|\nabla_{\theta_e} f(\bar{\theta}_1; \bar{\theta}_2)\| \leq C_{f_{\theta_e}}$  for some  $C_{f_{\theta_e}} > 0$ .

972 *Proof.*

$$973 \quad \|J_{ed}(\pi_{\theta_l, \theta_e})\| \leq E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t r_{ld} |s^0 = s| \right] \stackrel{(i)}{\leq} \frac{r_{ld}}{1-\gamma},$$

976 where (i) uses the close form of geometric series.

977 So the cumulative expert-dependent reward value  $J_{ed}$  is bounded, analogously, the cumulative  
978 expert-independent reward value  $J_{ei}$  is bounded.

$$\begin{aligned} 979 & \nabla_{\theta_l, \theta_e} f(\theta_l; \theta_e), \\ 980 & = \nabla_{\theta_e} (E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t [(\mu_l(s^t, a^t) - \mu_l(s^t)) (J_{ed}(s^t, a^t) + J_{ei}(s^t, a^t))^T] \right] + \mu_l(\pi_{\theta_l, \theta_e})), \\ 981 & \stackrel{(i)}{=} E^{\pi_{\theta_l, \theta_e}} \left[ \sum_{t=0}^{T-1} \gamma^t [(\nabla_{\theta_e} \mu_l(s^t, a^t) - \nabla_{\theta_e} \mu_l(s^t)) (J_{ed}(s^t, a^t) + J_{ei}(s^t, a^t))^T \right. \\ 982 & \quad \left. + (\mu_l(s^t, a^t) - \mu_l(s^t)) (\nabla_{\theta_e} J_{ed}(s^t, a^t) + \nabla_{\theta_e} J_{ei}(s^t, a^t))^T] \right] + \nabla_{\theta_e} \mu_l(\pi_{\theta_l, \theta_e}). \end{aligned}$$

988 Refer to the proof of Lemma 4, all elements in  $\nabla_{\theta_l, \theta_e} f(\theta_l; \theta_e)$  are finite, itself is bounded. Analo-  
989 gously,  $\nabla_{\theta_e} f(\theta_l; \theta_e)$  is bounded under the same proofing process.

990 Through the same procedure, we can prove  $\nabla_{\theta_e} f(\theta_l; \theta_e)$  is bounded.

991 All element of  $\nabla_{\theta_e} f(\theta_1; \theta_2)$  are bounded. Therefore, all elements for  $\nabla_{\theta_e} f(\theta_l; \theta_e)$  is bounded and  
992 it is bounded. □

#### 993 A.7 PROPERTIES OF THE APPROXIMATION ERRORS

994 Define  $b_{12}(k) = \hat{\nabla}_{\theta_l, \theta_e}^2 L(\theta_l(k), \theta_e(k)) - \nabla_{\theta_l, \theta_e}^2 L(\theta_l(k), \theta_e(k))$ ,  $b_{22}(k) = \hat{\nabla}_{\theta_e, \theta_e}^2 L(\theta_l(k), \theta_e(k)) -$   
995  $\nabla_{\theta_e, \theta_e}^2 L(\theta_l(k), \theta_e(k))$ ,  $b_1(k) = \hat{\nabla}_{\theta_l} f(\theta_l(k), \theta_e(k)) - \nabla_{\theta_l} f(\theta_l(k), \theta_e(k))$ ,  $b_2(k) =$   
996  $\hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k)) - \nabla_{\theta_e} f(\theta_l(k), \theta_e(k))$ ,  $b(k) = [\hat{\nabla}_{\theta_e, \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k)) -$   
997  $[\nabla_{\theta_e, \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \nabla_{\theta_e} f(\theta_l(k), \theta_e(k))$  at iteration  $k$ ,  $b_a(k) = \hat{\nabla} f(\theta_l(k), \theta_e(k)) -$   
998  $\nabla f(\theta_l(k), \theta_e(k))$  at the iteration  $k$ . Through the same procedure in the proof of the Lemma 5, we  
999 can get conclude that each element of forth-order gradient of  $L(\theta_l, \theta_e)$  is bounded by the constant  
1000  $C_L$  and each element of the third-order gradient of  $f(\theta_l(k), \theta_e(k))$  is bounded by  $C_f$ .

1001 **Lemma 6.** *The biases  $b_{12}(k)$ ,  $b_{22}(k)$ ,  $b_1(k)$ ,  $b_2(k)$ ,  $b(k)$ , and  $b_a(k)$  are bounded,*

$$\begin{aligned} 1002 & \|b_{12}(k)\| \leq \frac{C_L p^2(k) \sqrt{m}}{6} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}, \\ 1003 & \|b_{22}(k)\| \leq \frac{C_L p^2(k) \sqrt{m}}{6} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}, \\ 1004 & \|b_1(k)\| \leq \frac{C_f p^2(k) \sqrt{n}}{6} \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\}, \\ 1005 & \|b_2(k)\| \leq \frac{C_f p^2(k) \sqrt{m}}{6} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}, \\ 1006 & \|b(k)\| \leq \frac{2C_{f\theta_e} (C_L p^2(k) + C_f p^2(k))}{6\sqrt{m}\lambda^2} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}, \\ 1007 & \|b_a(k)\| \\ 1008 & \leq \frac{C_f p^2(k) \sqrt{n}}{6} \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\} \\ 1009 & + \frac{2C_{L\theta_l\theta_e} C_L p^2(k) + 2C_{L\theta_l\theta_e} \sqrt{m}\lambda C_f p^2(k)}{6\sqrt{m}\lambda^2} \\ 1010 & \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\} + \frac{C_{f\theta_e} C_L p^2(k)}{6m\lambda^3} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\} \\ 1011 & + \frac{2p^4(k) C_L (C_{f\theta_e} C_L + \sqrt{m}\lambda C_f)}{36\lambda^2} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}^2, \end{aligned}$$

1026 *Proof.* According to the Lemma 1 in (Spall, 1992), the approximation error  $b_{12l}(k)$  for  
 1027  $\hat{\nabla}_{\theta_l, \theta_e}^2 L(\theta_l, \theta_e)$  is :

$$1028 \quad b_{12l}(k) \leq \frac{C_{LP}^2(k)}{6} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\},$$

1031 where  $b_{12l}(k)$  represent the  $l$ th term of the bias  $b_{12l}(k)$  at  $k$ th iteration.

$$1032 \quad \|b_{12}(k)\| \leq \frac{C_{LP}^2(k) \sqrt{m}}{6} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}.$$

1033 Analogously, the  $b_{22}(k)$ ,  $b_1(k)$  and  $b_2(k)$  are also bounded.

1034 With the  $\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))$  and  $\hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k))$ , we estimate the  
 1035  $[\nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \nabla_{\theta_e} f(\theta_l(k), \theta_e(k))$  through conjugate gradient.

$$1036 \quad \|\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))\| = \|\nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k)) + b_{22}(k)\| \geq \|\lambda I + b_{22}(k)\|.$$

1037 Then we can tune the parameters of  $b_{22}(k)$  such as  $\Delta$ ,  $\alpha_0$  and  $\alpha_l$  to let  $\|[\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l, \theta_e)]^{-1}\| \leq \frac{2}{\sqrt{m\lambda}}$   
 1038 and make sure the  $\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l, \theta_e) \geq \frac{\lambda}{2} I$

$$1039 \quad \|b(k)\|,$$

$$1040 \quad = \|E[[\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k))$$

$$1041 \quad - [\nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \nabla_{\theta_e} f(\theta_l(k), \theta_e(k))]\|,$$

$$1042 \quad \stackrel{(iv)}{\leq} \|E[\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} (\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k)) - \nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k)))$$

$$1043 \quad [\nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1}]\| \|\nabla_{\theta_e} f(\theta_l(k), \theta_e(k))\| + \|\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} b_2(k)\|,$$

$$1044 \quad \stackrel{(v)}{\leq} \|[\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1}\| \|\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k)) - \nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))\|$$

$$1045 \quad [\nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \|\nabla_{\theta_e} f(\theta_l(k), \theta_e(k))\| + \|\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1}\| \|b_2(k)\|,$$

$$1046 \quad \stackrel{(vii)}{\leq} \frac{2C_{f\theta_e} \|b_{22}(k)\| + 2\sqrt{m\lambda} \|b_2(k)\|}{m\lambda^2},$$

$$1047 \quad \leq \frac{2C_{f\theta_e} (C_{LP}^2(k) + C_{fp}^2(k))}{6\sqrt{m\lambda^2}} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\},$$

1048 where (vii) uses the result of Lemma (5).

$$1049 \quad \|b_a(k)\| = E[\|\hat{\nabla} f(\theta_l(k), \theta_e(k))\| - \|\nabla f(\theta_l(k), \theta_e(k))\|,$$

$$1050 \quad \stackrel{(vii)}{=} E[\|\hat{\nabla}_{\theta_l} f(\theta_l(k), \theta_e(k)) - \hat{\nabla}_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k)) [\hat{\nabla}_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1}$$

$$1051 \quad \hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k))\| - \|\nabla f(\theta_l(k), \theta_e(k))\|,$$

$$1052 \quad = \|\nabla_{\theta_l} f(\theta_l(k), \theta_e(k)) + b_1(k) - (\nabla_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k)) + b_{12}(k))$$

$$1053 \quad ([\nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \nabla_{\theta_e} f(\theta_l(k), \theta_e(k)) + b(k))\| - \|\nabla f(\theta_l(k), \theta_e(k))\|$$

$$1054 \quad \stackrel{(v)}{\leq} \|b_1(k)\| + \|\nabla_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k))\| \|b(k)\|$$

$$1055 \quad + \|[\nabla_{\theta_e \theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1}\| \|\nabla_{\theta_e} f(\theta_l(k), \theta_e(k))\| \|b_{12}(k)\| + \|b_{12}(k)\| \|b(k)\|,$$

$$1056 \quad \stackrel{(vii)}{\leq} \frac{C_{fp}^2(k) \sqrt{n}}{6} \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\}$$

$$1057 \quad + \frac{2C_{L\theta_l \theta_e} C_{LP}^2(k) + 2C_{L\theta_l \theta_e} \sqrt{m\lambda} C_{fp}^2(k)}{6\sqrt{m\lambda^2}} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}$$

$$1058 \quad + \frac{C_{f\theta_e} C_{LP}^2(k)}{6m\lambda^3} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}$$

$$1059 \quad + \frac{2p^4(k) C_L (C_{f\theta_e} C_L + \sqrt{m\lambda} C_f)}{36\lambda^2} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\}^2,$$

where the first (vii) uses the result of Lemma (2), and the second (vii) uses the result of Lemma (5).  $\square$

### A.7.1 PROOF OF THEOREM 2

From the Lemma 2.2 of (Ghadimi & Wang, 2018),  $\|\nabla f(\theta_l(k), \theta_e(k)) - \nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\| \leq C\|\theta_e^*(\theta_l(k)) - \theta_e(k)\|$ ,  $\|\nabla f(\theta_l(k'), \theta_e^*(\theta_l(k'))) - \nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\| \leq L_f\|\theta_l(k') - \theta_l(k)\|$  where  $C = L_{f_{\theta_l}} + \frac{L_{f_{\theta_e}} C_{L_{\theta_l \theta_e}}}{\lambda} + C_{f_{\theta_e}} [\frac{L_{L_{\theta_l \theta_e}}}{\lambda} + \frac{L_{L_{\theta_e \theta_e}} C_{L_{\theta_l \theta_e}}}{\lambda^2}]$ ,  $L_f = \frac{(\bar{L}_{f_{\theta_e}} + C) C_{L_{\theta_l \theta_e}}}{\lambda} + L_{f_{\theta_l}} + C_{f_{\theta_e}} [\frac{\bar{L}_{L_{\theta_l \theta_e}} C_{f_{\theta_e}}}{\lambda} + \frac{\bar{L}_{L_{\theta_e \theta_e}} C_{L_{\theta_l \theta_e}}}{\lambda^2}]$ ,  $Q_L = \frac{L_{L_{\theta_e}}}{\lambda}$  denotes the condition number of  $L(\theta_l, \theta_e)$ ,  $M = \max_{\theta_l \in \theta_l} \|\theta_e(0) - \theta_e^*(\theta_l)\|$ ,  $D_{\theta_l} = \max_{x, y \in \theta_l} \{\|x - y\|\}$ . The proof of Theorem 2 is as follows:

*Proof.* First we compute the variance of  $\hat{\nabla} f(\theta_l(k), \theta_e(k))$ .

$$\begin{aligned}
& \|\hat{\nabla} f(\theta_l(k), \theta_e(k))\|, \\
& \stackrel{(vii)}{=} \|\hat{\nabla}_{\theta_l} f(\theta_l(k), \theta_e(k)) - \hat{\nabla}_{\theta_l}^2 L(\theta_l(k), \theta_e(k)) [\hat{\nabla}_{\theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \hat{\nabla}_{\theta_e} f(\theta_l(k), \theta_e(k))\|, \\
& = \|\nabla_{\theta_l} f(\theta_l(k), \theta_e(k)) + b_1(k) - (\nabla_{\theta_l}^2 L(\theta_l(k), \theta_e(k)) + b_{12}(k)) \\
& ([\nabla_{\theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1} \nabla_{\theta_e} f(\theta_l(k), \theta_e(k)) + b(k))\|, \\
& \stackrel{(v)}{\leq} \|\nabla_{\theta_l} f(\theta_l(k), \theta_e(k))\| + \|b_1(k)\| + \|\nabla_{\theta_l}^2 L(\theta_l(k), \theta_e(k))\| \|[\nabla_{\theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1}\| \\
& \|\nabla_{\theta_e} f(\theta_l(k), \theta_e(k))\| + \|\nabla_{\theta_l}^2 L(\theta_l(k), \theta_e(k))\| \|b(k)\| + \|b_1(k)\| \|[\nabla_{\theta_e}^2 L(\theta_l(k), \theta_e(k))]^{-1}\| \\
& \|\nabla_{\theta_e} f(\theta_l(k), \theta_e(k))\| + \|b_1(k)\| \|b(k)\|, \\
& \stackrel{(vii)}{\leq} C_{f_{\theta_l}} + \frac{C_f p^2(k) \sqrt{n}}{6} \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\} \\
& + \frac{C_{L_{\theta_l \theta_e}} C_{f_{\theta_e}}}{\sqrt{m} \lambda} + \frac{2C_{L_{\theta_l \theta_e}} C_{f_{\theta_e}} (C_L p^2(k) + C_f p^2(k))}{6\sqrt{m} \lambda^2} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\} \\
& + \frac{C_{f_{\theta_e}} C_f p^2(k) \sqrt{n}}{\sqrt{m} \lambda} \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\} \\
& + \frac{2C_f p^2(k) \sqrt{n} C_{f_{\theta_e}} (C_L p^2(k) + C_f p^2(k))}{36\sqrt{m} \lambda^2} \\
& \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\},
\end{aligned}$$

where the first (vii) uses the result of Lemma (2), and the second (vii) uses the result of Lemma (5) and Lemma (6).

Since  $\hat{\nabla} f(\theta_l(k), \theta_e(k))$  is bounded,

$$\begin{aligned}
& \text{Var}(\hat{\nabla} f(\theta_l(k), \theta_e(k))), \\
& \leq (C_{f_{\theta_l}} + \frac{C_f p^2(k) \sqrt{n}}{6} \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\} \\
& + \frac{C_{L_{\theta_l \theta_e}} C_{f_{\theta_e}}}{\sqrt{m} \lambda} + \frac{2C_{L_{\theta_l \theta_e}} C_{f_{\theta_e}} (C_L p^2(k) + C_f p^2(k))}{6\sqrt{m} \lambda^2} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\} \\
& + \frac{C_{f_{\theta_e}} C_f p^2(k) \sqrt{n}}{\sqrt{m} \lambda} \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\} \\
& + \frac{2C_f p^2(k) \sqrt{n} C_{f_{\theta_e}} (C_L p^2(k) + C_f p^2(k))}{36\sqrt{m} \lambda^2} \\
& \{[n^3 - (n-1)^3] \alpha_l^2 + (n-1)^3 \alpha_l \alpha_0^3\} \{[m^3 - (m-1)^3] \alpha_l^2 + (m-1)^3 \alpha_l \alpha_0^3\})^2.
\end{aligned}$$

Then we need to find the total bias, the total bias  $b_t(k)$  is the sum of the bias from approximation and  $\nabla f(\theta_l(k), \theta_e(k)) - \nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))$

1134  
1135  
1136  
1137  
1138  
1139  
1140

$$\begin{aligned} & \|b_t(k)\|, \\ & = \|E[\nabla f(\theta_l(k), \theta_e(k))] - \nabla f(\theta_l(k), \theta_e^*(k))\|, \\ & \stackrel{(vii)}{\leq} \|b_a(k)\| + \|r_{li}(\frac{Q_L-1}{Q_L+1})^{t_k} \|\theta_e(0) - \theta_e^*(\theta_l(k))\|, \end{aligned}$$

1141  
1142  
1143

where (vii) adds the approximation error  $r_{li}(\frac{Q_L-1}{Q_L+1})^{t_k} \|\theta_e(0) - \theta_e^*(\theta_l(k))\|$  from the lower level. Next step is to find the bound for  $E[\|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2]$ .

1144  
1145  
1146  
1147  
1148  
1149

$$\begin{aligned} & f(\theta_l(k+1), \theta_e^*(\theta_l(k+1))), \\ & \stackrel{(vi)}{\leq} f(\theta_l(k), \theta_e^*(\theta_l(k))) + \langle \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))), \theta_l(k+1) - \theta_l(k) \rangle + \frac{L_f}{2} \|\theta_l(k+1) - \theta_l(k)\|^2, \\ & = f(\theta_l(k), \theta_e^*(\theta_l(k))) - \alpha_k \langle \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))), \hat{\nabla} f(\theta_l(k), \theta_e(k)) \rangle + \frac{L_f \alpha_k^2}{2} \|\hat{\nabla} f(\theta_l(k), \theta_e(k))\|^2. \end{aligned}$$

1150  
1151

The expectation of  $f(\theta_l(k+1), \theta_e^*(\theta_l(k+1)))$  becomes:

1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166

$$\begin{aligned} & E[f(\theta_l(k+1), \theta_e^*(\theta_l(k+1)))], \\ & \leq f(\theta_l(k), \theta_e^*(\theta_l(k))) - \alpha_k \langle \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))), \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))) \rangle + b_t(k) \\ & + \frac{L_f \alpha_k^2}{2} E\|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k))) + \hat{\nabla} f(\theta_l(k), \theta_e(k)) - \nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2, \\ & \leq f(\theta_l(k), \theta_e^*(\theta_l(k))) - \alpha_k \langle \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))), \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))) \rangle + b_t(k) \\ & + \frac{L_f \alpha_k^2}{2} Var(\hat{\nabla} f(\theta_l(k), \theta_e(k))) + \frac{L_f \alpha_k^2}{2} E\|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2 \\ & + L_f \alpha_k^2 \langle \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))), b_t(k) \rangle, \\ & = f(\theta_l(k), \theta_e^*(\theta_l(k))) - (\alpha_k - \frac{L_f \alpha_k^2}{2}) \|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2 \\ & - (\alpha_k - L_f \alpha_k^2) \langle \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))), b_t(k) \rangle + \frac{L_f \alpha_k^2}{2} Var(\hat{\nabla} f(\theta_l(k), \theta_e(k))) + \frac{L_f \alpha_k^2}{2} \|b_t(k)\|^2. \end{aligned}$$

1167  
1168

Choose  $\alpha_k \leq \frac{1}{L_f}$  and with the fact  $2\langle \nabla f(\theta_l(k), \theta_e^*(\theta_l(k))), b_t(k) \rangle \leq \|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2 + \|b_t(k)\|^2$ .

1169  
1170  
1171  
1172  
1173  
1174  
1175

$$\begin{aligned} & E[f(\theta_l(k+1), \theta_e^*(\theta_l(k+1)))], \\ & \leq f(\theta_l(k), \theta_e^*(\theta_l(k))) - \frac{\alpha_k}{2} \|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2 \\ & + \frac{\alpha_k}{2} \|b_t(k)\|^2 + \frac{L_f \alpha_k^2}{2} Var(\hat{\nabla} f(\theta_l(k), \theta_e(k))), \end{aligned}$$

1176  
1177

Rearrange terms,

1178  
1179  
1180  
1181  
1182  
1183

$$\begin{aligned} & \sum_{k=0}^{K-1} \frac{\alpha_k}{2} E[\|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2], \\ & \leq f(\theta_l(0), \theta_e^*(\theta_l(0))) - f^* + \sum_{k=0}^{K-1} (\frac{\alpha_k}{2} \|b_t(k)\|^2 + \frac{L_f \alpha_k^2}{2} Var(\hat{\nabla} f(\theta_l(k), \theta_e(k)))). \end{aligned}$$

1184  
1185  
1186  
1187

For  $\|b_t(k)\|^2$ , it is the linear combination of  $p^4(k), p^6(k), p^8(k), p^2(k)(\frac{Q_L-1}{Q_L+1})^{t_k}, p^4(k)(\frac{Q_L-1}{Q_L+1})^{t_k}, (\frac{Q_L-1}{Q_L+1})^{2t_k}$ . For  $Var(\hat{\nabla} f(\theta_l(k), \theta_e(k)))$ , it is the linear combination of  $1, p^2(k), p^4(k), p^6(k), p^8(k)$ . For simplification, we use  $C_{si} > 0, i = 1, 2, \dots$  to represent the constant for all combinations of terms involve  $p(k), \alpha_k$  and  $t_k$ . Then we can continue the calculation as follows:

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

$$\begin{aligned}
& \frac{1}{K} \sum_{k=0}^{K-1} E[\|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2], \\
& \leq \frac{2}{K\alpha_k} f(\theta_l(0), \theta_e^*(\theta_l(0))) - f^* + \frac{1}{K} \sum_{k=0}^{K-1} (\|b_t(k)\|^2 + L_f \alpha_k \text{Var}(\hat{\nabla} f(\theta_l(k), \theta_e(k)))), \\
& \stackrel{(viii)}{\leq} \frac{2(f(\theta_l(0), \theta_e^*(\theta_l(0))) - f^*)}{\alpha_k K} + \frac{C_{s1}}{K} \sum_{k=0}^{K-1} p^4(k) + \frac{C_{s2}}{K} \sum_{k=0}^{K-1} p^6(k) + \frac{C_{s3}}{K} \sum_{k=0}^{K-1} p^8(k) \\
& + \frac{C_{s4}}{K} \sum_{k=0}^{K-1} p^2(k) \left(\frac{Q_L - 1}{Q_L + 1}\right)^{t_k} + \frac{C_{s5}}{K} \sum_{k=0}^{K-1} p^4(k) \left(\frac{Q_L - 1}{Q_L + 1}\right)^{t_k} + \frac{C_{s6}}{K} \sum_{k=0}^{K-1} \left(\frac{Q_L - 1}{Q_L + 1}\right)^{2t_k} \\
& + \frac{C_{s7}}{K} \sum_{k=0}^{K-1} \alpha_k + \frac{C_{s8}}{K} \sum_{k=0}^{K-1} \alpha_k p^2(k) + \frac{C_{s9}}{K} \sum_{k=0}^{K-1} \alpha_k p^4(k) + \frac{C_{s10}}{K} \sum_{k=0}^{K-1} \alpha_k p^6(k) \\
& + \frac{C_{s11}}{K} \sum_{k=0}^{K-1} \alpha_k p^8(k),
\end{aligned}$$

where (viii) expands each term of  $\|b_t(k)\|^2$  and  $\text{Var}(\hat{\nabla} f(\theta_l(k), \theta_e(k)))$ . Choose  $p(k) = \frac{1}{k}$ ,  $\alpha_k = \frac{1}{L_f \sqrt{K}}$ ,  $t_k = \lceil \frac{\sqrt{k+1}}{2} \rceil$ . Since  $0 \leq \frac{Q_L - 1}{Q_L + 1} < 1$ , we can conclude  $\sum_{k=0}^{K-1} p^2(k) \left(\frac{Q_L - 1}{Q_L + 1}\right)^{t_k} < \sum_{k=0}^{K-1} p^2(k)$  when  $\frac{Q_L - 1}{Q_L + 1} \neq 0$ , then the convergence rate is as follows:

$$\frac{1}{K} \sum_{k=0}^{K-1} E[\|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2] \leq \frac{C_{s12}}{\sqrt{K}} + \frac{C_{s13}}{K} + \frac{C_{s14}}{K\sqrt{K}}.$$

As  $K \rightarrow \infty$ ,  $\frac{1}{K} \sum_{k=0}^{K-1} E[\|\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))\|^2] \rightarrow 0$ , which shows that  $E[\nabla f(\theta_l(k), \theta_e^*(\theta_l(k)))]$  decreases at the rate of  $\mathcal{O}(\frac{1}{\sqrt{K}} + \frac{1}{K} + \frac{1}{K\sqrt{K}})$ .  $\square$

#### A.7.2 PROOF OF COROLLARY 2.1

Define the cumulative reward function of the expert as  $J_e(\pi) \triangleq E^\pi[\sum_{t=0}^{T-1} \gamma^t r_{\theta_e}(s^t, a^t)]$ . If  $r_{\theta_e}$  is a linear reward function, we have  $r_{\theta_e} \triangleq \langle \theta_e, \phi(s, a_e) \rangle$  where the feature  $\phi(s, a_e) \in \mathbb{R}^{d_{\theta_e}}$  and  $d_{\theta_e}$  is the dimension of  $\theta_e$ . Then the expert's feature expectation is formulated as  $\mu_f(\pi) \triangleq E^\pi[\sum_{t=0}^{T-1} \gamma^t \phi(s^t, a_e^t)]$ . From Theorem 2, we can get  $\|\mu_f(\pi_{\theta_e}) - \mu_f(\pi_e)\|^2$  decreases in  $\mathcal{O}(\frac{1}{\sqrt{K}})$ .

*Proof.*

$$\begin{aligned}
J_e(\pi_{\theta_e}) - J_e(\pi_e) &= \langle \theta_e, \mu_f(\pi_{\theta_e}) - \mu_f(\pi_e) \rangle \\
&\leq \max_{\theta_e \in \Theta} \langle \theta_e, \mu_f(\pi_{\theta_e}) - \mu_f(\pi_e) \rangle \\
&\stackrel{(v)}{\leq} \max_{\theta_e \in \Theta} \|\theta_e\| \|\mu_f(\pi_{\theta_e}) - \mu_f(\pi_e)\|, \\
&= \|\mu_f(\pi_{\theta_e}) - \mu_f(\pi_e)\|, \\
&\stackrel{(vii)}{\leq} \frac{C_{15}}{\sqrt{K}},
\end{aligned}$$

where (vii) uses the result  $\|\mu_f(\pi_{\theta_e}) - \mu_f(\pi_e)\|^2$  decreases in  $\mathcal{O}(\frac{1}{\sqrt{K}})$  and  $C_{15}$  is the constant number which includes all influence factors other than  $K$ .  $\square$

#### A.8 PROOF OF THEOREM 1

The prove of Theorem 1 is as follows: According to the Algorithm 1 in (Ziebart et al., 2008), we need to compute the state frequency for the  $\mu_e(\pi_{\theta_l, \theta_e})$ . For each state-action pair, it needs to recursively compute for up to  $T$  iterations. As there are  $T$  state-action

1242 pairs in one demonstration, the computational complexity for the  $\mu_e(\pi_{\theta_l, \theta_e})$  is  $\mathcal{O}(T^2)$   
 1243 as we see  $T$  as the deterministic factor. Analogously, the computational complexity for  
 1244  $\mu_l(\pi_{\theta_l, \theta_e}), \mu_e(s, a), \mu_e(s), \mu_l(s, a), \mu_l(s), J_{ei}(s, a), J_{ed}(\pi_{\theta_l, \theta_e}), J_{ei}(\pi_{\theta_l, \theta_e})$  are  $\mathcal{O}(T^2)$ . We can  
 1245 see that these factors with computational complexity  $\mathcal{O}(T^2)$  are expected to be calculated through  
 1246 backpropagation, and the backpropagation leads to  $\mathcal{O}(T^2)$ . However, in coding, we can calculate  
 1247 these factors by sampling a finite number of trajectories, and the computational complexity becomes  
 1248  $\mathcal{O}(T)$ .

1249  
 1250 For SPSA, the required terms are  $f(\theta_l, \theta_e), \nabla_{\theta_l} L(\theta_l, \theta_e), \nabla_{\theta_e} L(\theta_l, \theta_e)$ , these terms are the linear  
 1251 combination of  $\mathcal{O}(T)$  computational complexity terms, so their computational complexities are also  
 1252  $\mathcal{O}(T)$ .

1253  
 1254 If we directly compute  $\nabla_{\theta_l} f(\theta_l(k), \theta_e(k)), \nabla_{\theta_l \theta_e}^2 L(\theta_l(k), \theta_e(k)), \nabla_{\theta_e}^2 L(\theta_l(k), \theta_e(k)),$   
 1255 and  $\nabla_{\theta_e} f(\theta_l(k), \theta_e(k))$  instead of approximating, use the expression of  $\nabla_{\theta_l} f(\theta_l, \theta_e) =$   
 1256  $E^{\pi_{\theta_l, \theta_e}} [\sum_{t=0}^{T-1} \gamma^t [(\mu_l(s^t, a^t) - \mu_l(s^t))(J_{ed}(s^t, a^t) + J_{ei}(s^t, a^t))^T]] + \mu_l(\pi_{\theta_l, \theta_e})$  as an exam-  
 1257 ple,  $\mu_l(s)$  is inside an expectation from  $t = 0$  to  $T - 1$ , we need to sum up  $\mu_l(s)$  for  $T$  times.  
 1258 As a result, the computational complexity for  $\nabla_{\theta_l} f(\theta_l, \theta_e)$  is  $\mathcal{O}(T^2)$ . Analogously,  $\nabla_{\theta_e} f(\theta_l, \theta_e),$   
 1259  $\nabla_{\theta_l \theta_e}^2 L(\theta_l, \theta_e), \nabla_{\theta_e \theta_e}^2 L(\theta_l, \theta_e)$  are all same.

1260  
 1261 Back to SPSA, more policies need to be found compared to directly compute the hypergradient. Use  
 1262 soft q learning (Haarnoja et al., 2017) as an example, RL and MARL algorithms can be considered  
 1263 as sampling a finite number trajectories for each epoch. Therefore the computational cost for each  
 1264 epoch is  $\mathcal{O}(T)$  and the overall computational cost is  $\mathcal{O}(eT)$  where  $e$  is the total number of epochs.  
 1265 As we consider  $T$  as the decision factor, the computational cost of the multi-agent RL is  $\mathcal{O}(T)$ ,  
 1266 which matches the computational complexity of SPSA.

1267  
 1268 As a result, the computational complexity of SPSA is dominated by  $\mathcal{O}(T)$ , and directly calculating  
 1269 the hypergradient is dominated by  $\mathcal{O}(T^2)$ .

## 1270 A.9 EXPERIMENT DETAIL

1271  
 1272 The details of the experiments are shown in this section. All Python3 codes are run on a Windows  
 1273 10 desktop with 13th Gen Intel(R) Core(TM) i7-13700KF CPU and 32 GB of RAM. For each  
 1274 combination of algorithms and environments, we run 10 times to calculate mean values and standard  
 1275 deviations at each iteration. Then the calculated mean values and standard deviations are plotted as  
 1276 shown in Section 7 figures.

### 1277 A.9.1 MPE

1278  
 1279 The state, action, and observation spaces for the adversary and good agents are continuous. For the  
 1280 adversary, it can observe the relative distance to the landmarks and the good agents, therefore the  
 1281 observation of the adversary is  $o_a = [p_{l1} - p_a, p_{l2} - p_a, p_{g1} - p_a, p_{g2} - p_a]$  where  $p_{l1}$  is the position  
 1282 of the first landmark,  $p_{l2}$  is the position of the second landmark,  $p_{g1}$  is the position of the first good  
 1283 agent,  $p_{g2}$  is the position of the second good agent. For each good agent, it can observe the relative  
 1284 distance to the target landmark, the landmarks, the adversary, and another good agent, therefore the  
 1285 observations for two good agents are  $o_{g1} = [p_{tl} - p_{g1}, p_{l1} - p_{g1}, p_{l2} - p_{g1}, p_a - p_{g1}, p_{g2} - p_{g1}]$   
 1286 and  $o_{g2} = [p_{tl} - p_{g2}, p_{l1} - p_{g2}, p_{l2} - p_{g2}, p_a - p_{g2}, p_{g1} - p_{g2}]$  where  $p_{tl}$  is the position of the  
 1287 target landmark. The actions of the adversary and the good agents are the velocities between 0  
 1288 and 1 in four directions (left, right, down, up). Two good agents share the same return, which is  
 1289 rewarded based on the minimum distance of any agent to the target landmark and is penalized based  
 1290 on the distance between the adversary and the target landmark, therefore the reward of good agents  
 1291 is  $r_g = -\min(\|p_{tl} - p_{g1}\|_2, \|p_{tl} - p_{g2}\|_2) + \|p_{tl} - p_a\|_2$ . The reward of the adversary is based on  
 1292 the distance to the target of the adversary, therefore  $r_a = -\|p_{ta} - p_a\|_2$ , where  $p_{ta}$  is the position of  
 1293 the adversary's target. In our simulation, we consider observations as states of the MG, the distance  
 1294  $-\min(\|p_{tl} - p_{g1}\|_2, \|p_{tl} - p_{g2}\|_2)$  as the adversary-independent reward function, and the distance  
 1295  $\|p_{tl} - p_a\|_2$  as the feedback received by the good agents.

## A.9.2 SMAC

In this scenario, we can only access the state, observation, and action information of agents. The enemy is controlled by a built-in game AI. We consider observations as the states of the agents during learning. For each agent, it can observe the following information corresponding to another agent and enemy: relative distance, relative x, relative y, health, shield, and unit type. If an agent or an enemy is under attack, the shield reduces first and health reduces after the shield disappears. There are 7 possible actions for each agent: move north, move south, move east, move west, attack the enemy, stop, and no-op. At each time step, each agent can know which action among these 7 possible actions are available and the agent chooses one action from available actions. For example, the attack action is available when the enemy is in the shooting range of the agent. The agent can only choose no-op when its own health is 0. The agent gains rewards based on the damage dealt to the enemy and if the enemy is defeated.

## A.9.3 HUMAN-ROBOT INTERACTION

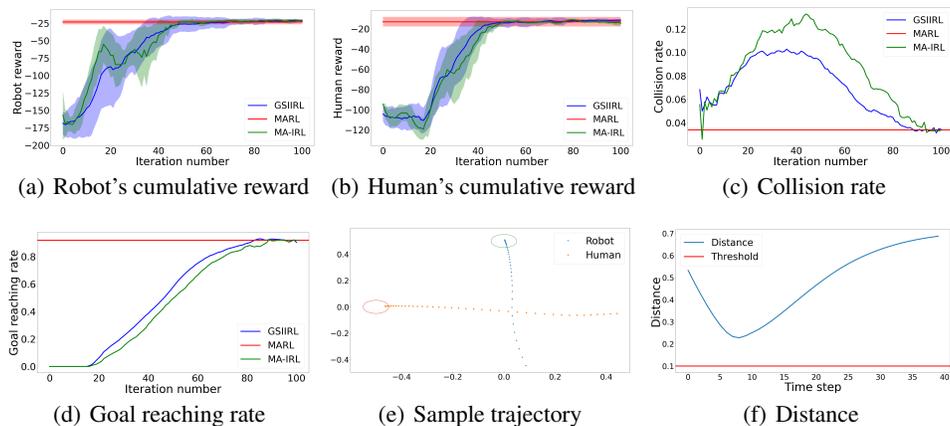


Figure 4: Human-robot interaction simulation results. The cumulative rewards are calculated as described in Figure 2. Similarly, the collision and goal-reaching rates are computed based on the policies corresponding to the learned reward functions. A collision is defined to occur when the distance between the human and the robot falls below a threshold. The goal is considered reached only when both the human and the robot reach their respective destinations.

Consider the human-robot interaction scenario described in the introduction, both the robot (learner) and the human (expert) aim to reach their destinations as quickly as possible. Additionally, the robot strives to maintain the human’s right-of-way. Both agents operate in continuous state and action spaces: each has a 4-dimensional state space and a 2-dimensional action space.

Figures 4(a) and 4(b) show that the robot’s and the human’s cumulative rewards converge to the ground-truth values, consistent with previous experiments. Figure 4(c) shows that the collision rate converges to its ground-truth value, indicating that the robot successfully learns a human-dependent reward function. Initially, with randomly generated reward functions, the robot and human move randomly, and collisions are rare. As the reward functions are iteratively updated, both agents move toward their goal locations, causing collisions and an increase in the collision rate. However, as training continues and the reward functions improve to prioritize collision avoidance, the collision rate decreases. Figure 4(d) shows that the goal-reaching rate also converges to the ground-truth value. With the learned reward functions, both agents reliably reach their goals. After 100 iterations of reward function updates, the learned robot reward function yields a high goal-reaching probability and a low collision probability. Figure 4(e) presents a sample trajectory under the learned reward functions, and Figure 4(f) plots the distance between the robot and the human along this trajectory. Both agents reach their goals while maintaining a separation above the required safety distance.

In the simulation, the policies for the human and the robot are calculated through Multi-Agent Deep Deterministic Policy Gradient (MADDPG) (Lowe et al., 2017). During the training process, noises are added to the action to increase the exploration. Once the policies are generated, further calculations use deterministic policies.

The state of the robot is its location  $s_r = (x_r, y_r) \in \mathbb{R}^2$ , the action of the robot includes the horizontal and vertical velocities and defined as  $a_r = (v_{rx}, v_{ry})$ , where  $v_{rx} \in [-0.1, 0.1]$ ,  $v_{ry} \in [-0.1, 0.1]$ . Similarly, the state of the human is  $s_h = (x_h, y_h) \in \mathbb{R}^2$ , the action of the human is  $a_h = (v_{hx}, v_{hy})$ , where  $v_{hx} \in [-0.1, 0.1]$ ,  $v_{hy} \in [-0.1, 0.1]$ . At each time step  $t$ , the robot chooses the action  $a_r(t)$  based on the current joint state  $(s_r(t), s_h(t))$  and moves to the next state  $x_r(t+1) = x_r(t) + v_{rx}(t)$ ,  $y_r(t+1) = y_r(t) + v_{ry}(t)$ . Analogously, the motion dynamics of the human is given by  $x_h(t+1) = x_h(t) + v_{hx}(t)$ ,  $y_h(t+1) = y_h(t) + v_{hy}(t)$ . In the experiment, the robot starts from an initial state  $s_r(0) \in [-0.25, 0.25] \times [-0.4, 0.5]$  and aims to reach a circle goal region whose center is at  $(0, 0.5)$  with radius 0.05. The human starts from  $s_h(0) \in [0.4, 0.5] \times [-0.25, 0.25]$  and aims to reach a circle goal region whose center is at  $(-0.5, 0)$  with radius 0.05. Both the robot and the human are penalized when a collision happens.

#### A.9.4 SECURITY

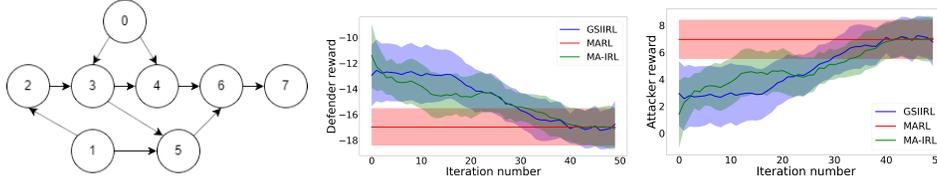


Figure 5: Cyber security simulation results. **Left:** Attack graph. **Middle:** Defender’s reward. **Right:** Attacker’s reward. The cumulative rewards are calculated in the same way described in Figure 2.

We conduct a cybersecurity experiment to evaluate the proposed algorithm. The experimental setup involves a defender (learner) and an attacker (expert) interacting on the attack graph in Figure 5. The attacker attempts to compromise nodes in the graph. The attacker’s objective is to compromise as many nodes as possible, whereas the defender’s objective is to protect the network by minimizing the number of compromised nodes. Both the learner and the expert have discrete state and action spaces, with cardinalities of 256 and 8, respectively.

Figure 5 shows that the cumulative rewards of the proposed algorithm converge to those of MARL, consistent with previous experiments.

There are 8 nodes and 10 edges. Each node represents a machine, and each edge represents an exploit between two nodes. The decision-making of the defender and the attacker is modeled as an MG. The state  $s \in \{0, 1\}^8$ , represents the condition of each node where the value 1 means the current node is compromised by the attacker, and the value 0 is vice versa. In each action pair, the attacker chooses one edge to attack, and the defender chooses one edge to block. Suppose the attacker chooses to attack the edge  $\{i, j\}$ . If the node  $i$  is already compromised and the defender does not block this edge, there is a probability for node  $j$  to be compromised. For other situations, the node  $j$  keeps clean. Each edge has a cost for the attacker to utilize and a cost for the defender to block. The attacker receives a reward when it successfully compromises a new node. The net reward of the attacker for each state-action pair is the sum of the reward and the cost. For the defender, the expert-dependent reward is the opposite of the attacker’s reward, and the expert-independent reward is the cost to block edges.

For the security simulation, the attack graph is randomly generated. We use Q-learning to find the policies for the attacker and the defender. During the training process, the attacker and the defender have a 70% possibility to choose between the best action with 60% possibility and the second best action with 40% possibility. Otherwise, the attacker and the defender randomly choose one action from the action space. When we exploit the learned policies, the attacker and the defender choose between the best action with 60% possibility and the second best action with 40% possibility.

#### A.9.5 CUMULATIVE REWARDS COMPARE TABLE

#### A.10 THE USE OF LARGE LANGUAGE MODELS (LLMs)

We confirm that LLM (ChatGPT 5) assistance was limited to improving grammar and readability.

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

Table 1: The cumulative reward comparison. **Top:** The cumulative reward of the learner. **Bottom:** The cumulative reward of the expert. MARL uses ground-truth reward functions. MA-IRL and GSIIRL use learned reward functions from the last iteration.

LEARNER	MARL	MA-IRL	GSIIRL
MPE	$5.03 \pm 1.66$	$5.12 \pm 4.08$	$4.84 \pm 1.81$
SMAC	$19.01 \pm 0.13$	$18.89 \pm 0.20$	$18.87 \pm 0.27$
CS	$-16.96 \pm 1.45$	$-17.08 \pm 3.21$	$-17.38 \pm 3.59$
HRI	$-20.22 \pm 3.34$	$-20.63 \pm 1.07$	$-21.09 \pm 1.87$
EXPERT	MARL	MA-IRL	GSIIRL
MPE	$-18.91 \pm 2.01$	$-20.37 \pm 4.27$	$-19.53 \pm 2.50$
SMAC	$19.01 \pm 0.13$	$18.89 \pm 0.20$	$18.87 \pm 0.27$
CS	$6.96 \pm 1.45$	$7.03 \pm 3.21$	$6.76 \pm 3.59$
HRI	$-12.98 \pm 4.78$	$-13.52 \pm 2.57$	$-13.16 \pm 1.03$