

AD HOC TEAMWORK & ITS CHALLENGES

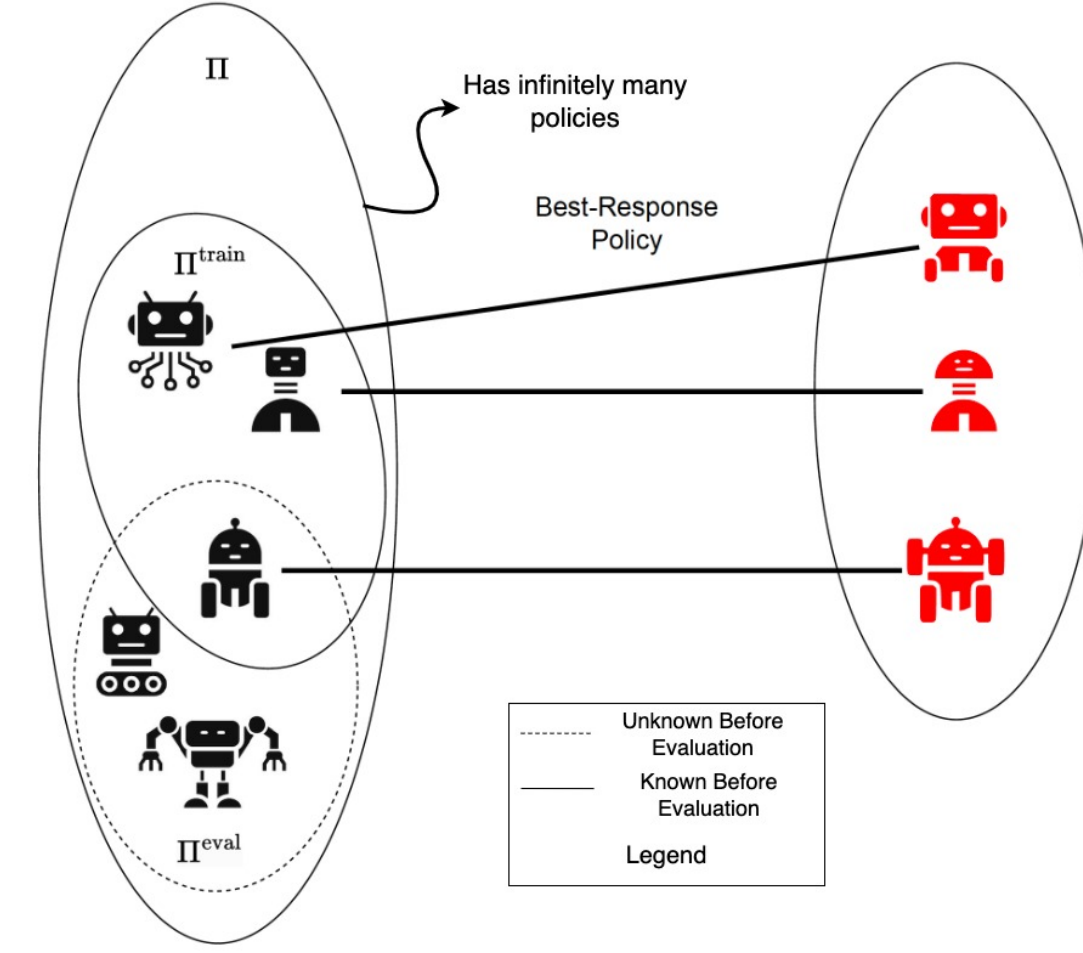
Ad Hoc Teamwork

- Create an adaptive agent (**learner**) that can collaborate with others without prior coordination mechanisms
- Train learner policy ($\pi^{*,i}$) with training set of teammate policies (Π^{train})
- Given a holdout set of teammate policies, Π^{eval} , evaluate the expected returns of $\pi^{*,i}$ defined below:

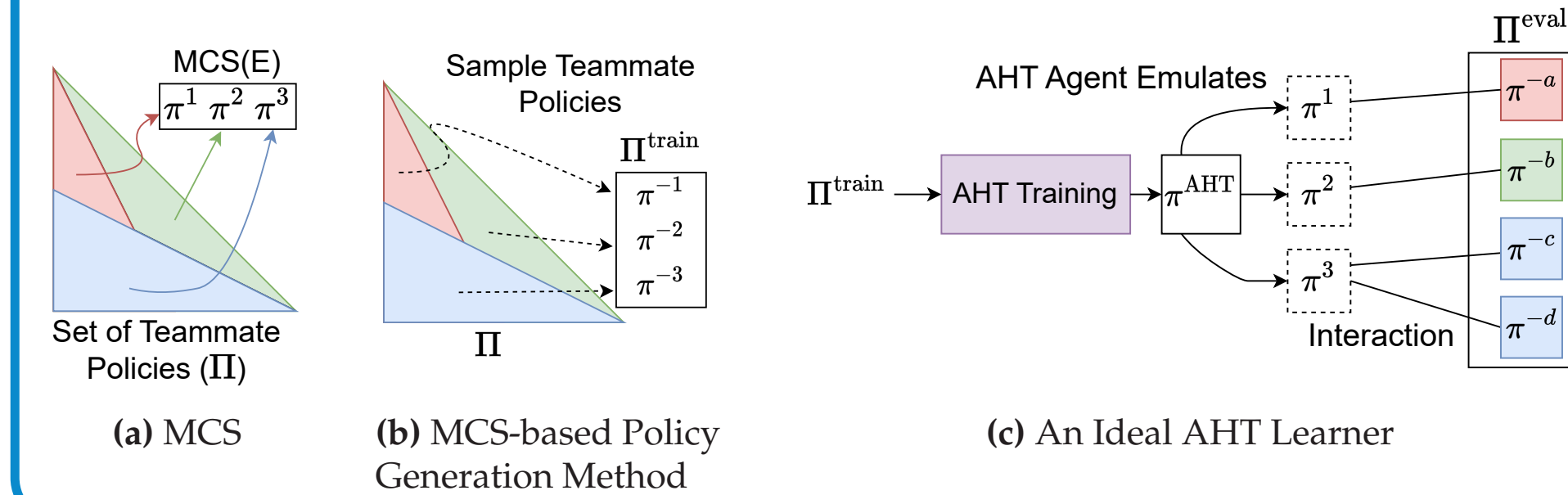
$$\mathbb{E}_{\pi^{-i} \sim \mathcal{U}(\Pi^{\text{eval}}), a_t^i \sim \pi^{*,i}, a_t^{-i} \sim \pi^{-i}, P, O} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right] \quad (1)$$

Challenges

- Π^{eval} is unknown when training the learner
- Cannot train with the infinite set of all teammate policies (Π)
- Π^{train} not necessarily representative of Π^{eval}



IDEAL TEAMMATE POLICY GENERATION METHOD



1. Find the minimum coverage set (MCS).
2. For each $\pi^i \in \text{MCS}$, sample π^{-i} that has π^i as its best-response (BR) policy and include as part of Π^{train} .

Goal:

- Enable learner to emulate BR policy to any $\pi^{-i} \in \Pi$.

LAGRANGIAN BEST RESPONSE DIVERSITY (L-BRDIV)

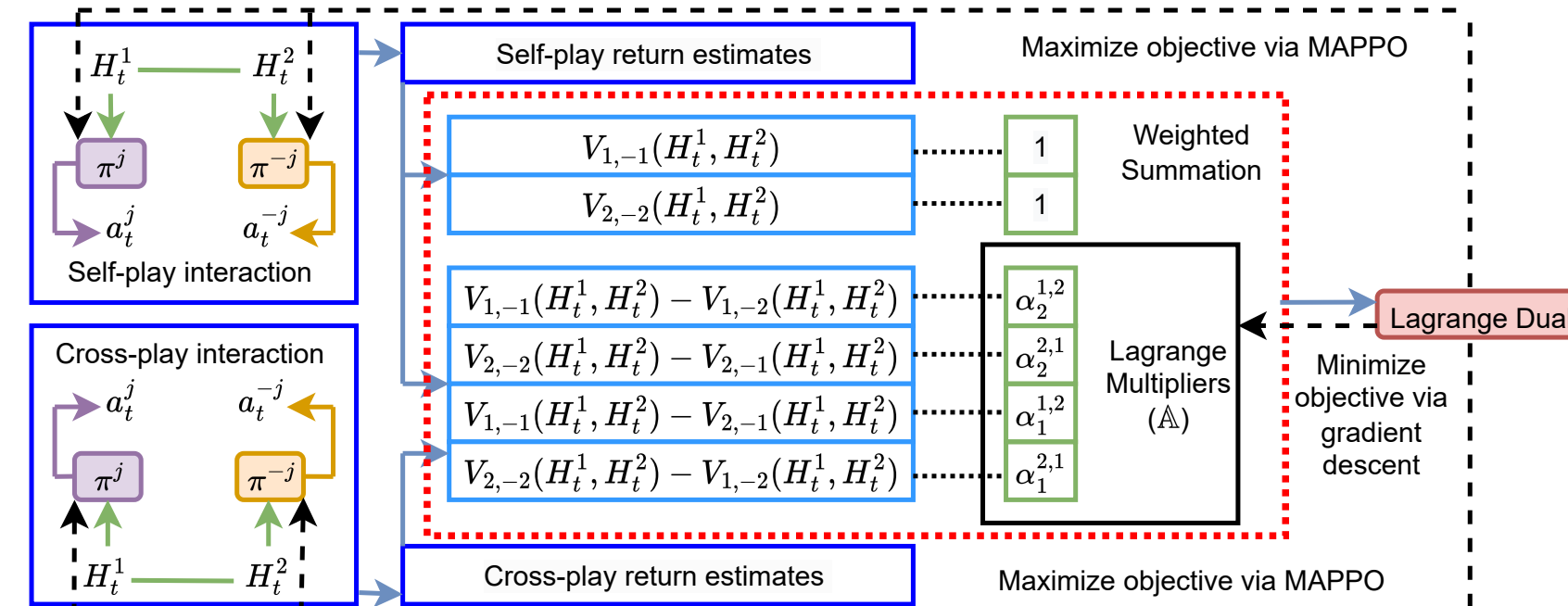
Generate $\Pi^{\text{train}} = \{\pi^{-i}\}_{i=1}^K$ and their set of BR policies $\{\pi^i\}_{i=1}^K$ by optimizing the following objective:

$$\max_{\substack{\{\pi^i\}_{i=1}^K \subseteq \Pi, \\ \{\pi^{-i}\}_{i=1}^K \subseteq \Pi}} \sum_{i \in \{1, 2, \dots, K\}} \mathbb{E}_{s \sim p_0} [\mathbf{R}_{i,-i}(H_t)], \quad (2)$$

with the following constraints that must be fulfilled for all $i, j \in \{1, 2, \dots, K\}$ and $i \neq j$:

$$\mathbb{E}_{s \sim p_0} [\mathbf{R}_{j,-i}(H_t)] + \tau \leq \mathbb{E}_{s \sim p_0} [\mathbf{R}_{i,-i}(H_t)], \quad (3)$$

$$\mathbb{E}_{s \sim p_0} [\mathbf{R}_{i,-j}(H_t)] + \tau \leq \mathbb{E}_{s \sim p_0} [\mathbf{R}_{i,-i}(H_t)]. \quad (4)$$

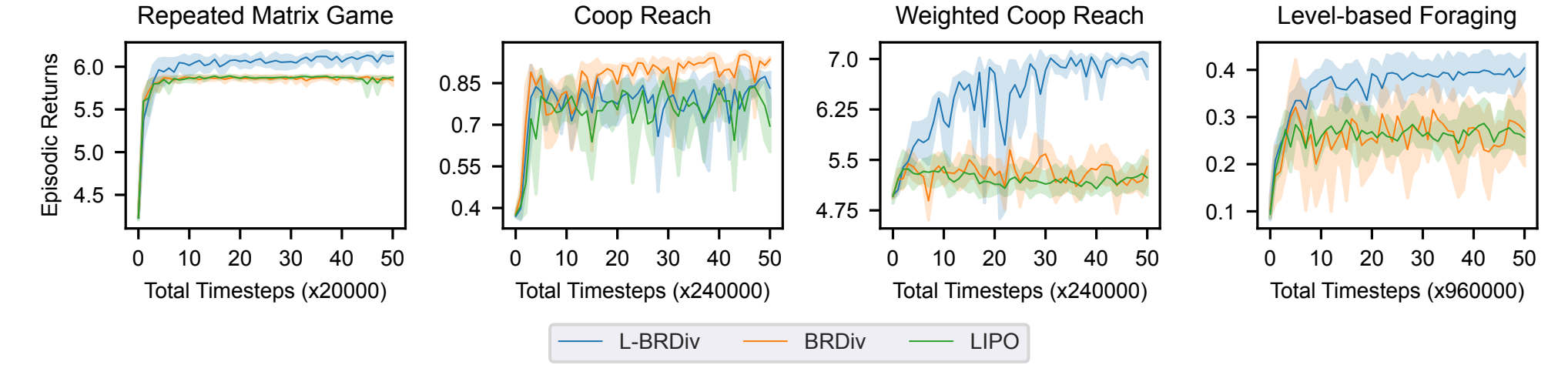


EXPERIMENT RESULTS

Generalization Experiments

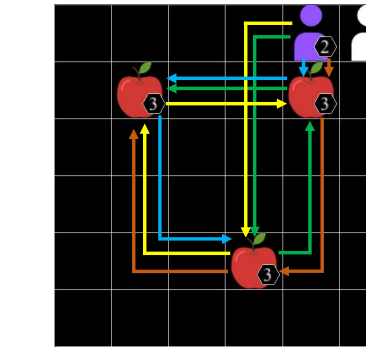
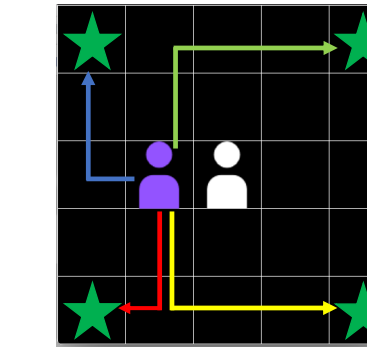
(Baselines) Teammate policy generation methods maximizing adversarial diversity

(Experiment Protocol) With Π^{train} generated by each compared method, train the learner and evaluate its returns when dealing with holdout policies in Π^{eval}



Analysis of Generated Policies

	$\pi(A)$	$\pi(B)$	$\pi(C)$
1	1	0	0
2	0	1	0
3	0	0	1



(a) Repeated Matrix Game

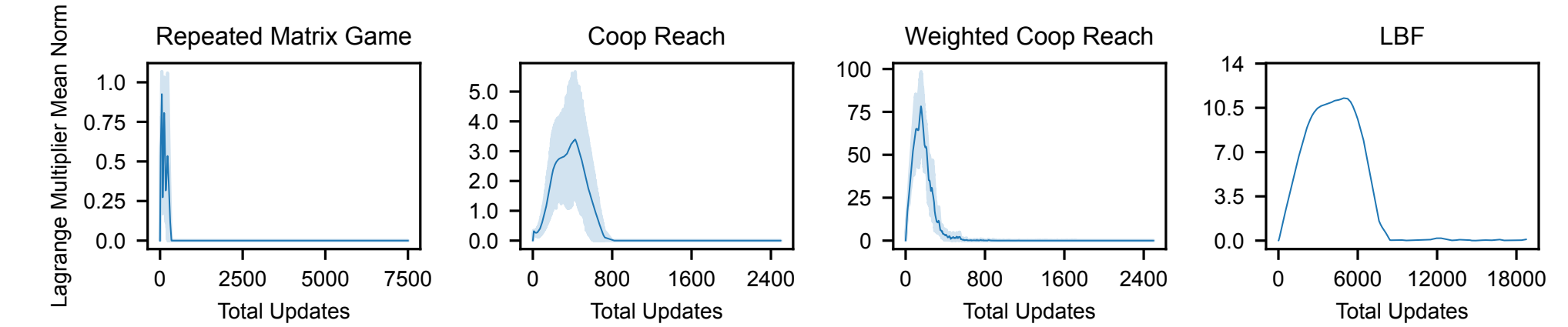
(b) (Weighted) Coop Reaching

(c) LBF

- L-BRDIV facilitates the discovery of Π^{train} having more members with different BR policies
- The BR policies of Π^{train} encompass all members of the MCS

Lagrange Multiplier Analysis

- Lagrange multipliers keep increasing while constraints are violated
- Eventually, the Lagrange multipliers converge to zero once constraints are fulfilled



SUMMARY & FUTURE WORK

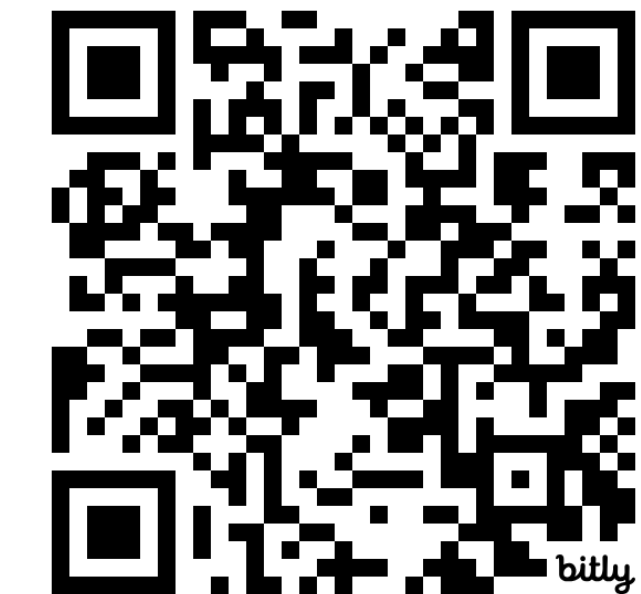
Our Contributions

- Important concept (i.e. minimum coverage sets) for generating Π^{train} that enable training robust AHT agents
- A Lagrangian multiplier-based teammate generation method (i.e. L-BRDIV) that outperforms existing state-of-the-art baselines

Future Work

- Extend agent generation method to general-sum games
- Generalizing to N-Player games

EXTERNAL LINKS



(a) Paper



(b) Code Repository