
ROTATE: Regret-driven Open-ended Training for Ad Hoc Teamwork

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning to collaborate with previously unseen partners is a fundamental general-
2 ization challenge in multi-agent learning, known as Ad Hoc Teamwork (AHT).
3 Existing AHT approaches often adopt a two-stage pipeline, where first, a fixed
4 population of teammates is generated with the idea that they should be repre-
5 sentative of the teammates that will be seen at deployment time, and second, an
6 AHT agent is trained to collaborate well with agents in the population. To date,
7 the research community has focused on designing separate algorithms for each
8 stage. This separation has led to algorithms that generate teammates with limited
9 coverage of possible behaviors, and that ignore whether the generated teammates
10 are easy to learn from for the AHT agent. Furthermore, algorithms for training
11 AHT agents typically treat the set of training teammates as static, thus attempting
12 to generalize to previously unseen partner agents without assuming any control
13 over the set of training teammates. This paper presents a unified framework for
14 AHT by reformulating the problem as an open-ended learning process between
15 an AHT agent and an adversarial teammate generator. We introduce ROTATE, a
16 regret-driven, open-ended training algorithm that alternates between improving
17 the AHT agent and generating teammates that probe its deficiencies. Experiments
18 across diverse two-player environments demonstrate that ROTATE significantly
19 outperforms baselines at generalizing to an unseen set of evaluation teammates,
20 thus establishing a new standard for robust and generalizable teamwork.

21 1 Introduction

22 As AI agents are deployed in diverse applications, it is increasingly crucial that they can collaborate
23 effectively with previously unseen AI agents and humans. While methods for training teams of
24 agents have been explored in cooperative multi-agent reinforcement learning (CMARL) [20, 44],
25 prior work highlighted that CMARL agents fail to perform optimally when collaborating with
26 unfamiliar teammates [54, 41]. Rather than learning strategies that are only effective against jointly
27 trained teammates, dealing with previously unseen teammates requires adaptive AI agents that
28 efficiently approximate the optimal strategy for collaborating with diverse teammates. The training
29 of such adaptive agents has been explored within ad hoc teamwork (AHT) [7, 51, 35] and zero-shot
30 coordination (ZSC) [22, 15, 34].

31 Most work has decomposed AHT learning into two stages [35], consisting of first creating a fixed
32 set of teammates, and then training an AHT agent using reinforcement learning (RL), based on
33 interactions with teammates sampled from the set. Methods that focus on AHT agent learning
34 typically rely on a human-designed heuristic-based or pretrained teammates [40, 66, 41] and therefore
35 struggle to handle novel behaviors outside the predefined set of teammates [52, 9]. Recent work
36 enhances the generalization capabilities of AHT agent learning methods by substituting the predefined
37 set of teammates with a generated collection of diverse teammates [34, 43], which are trained to

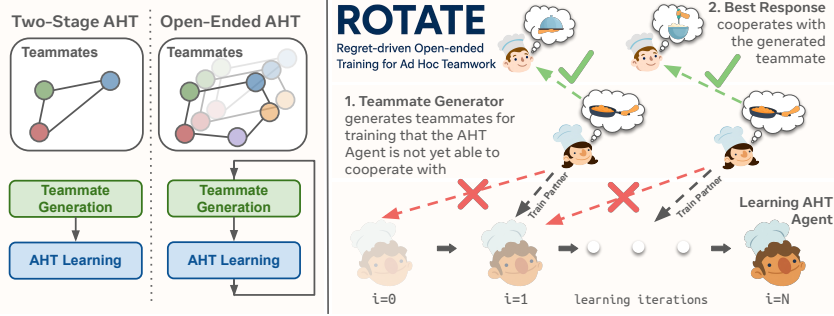


Figure 1: ROTATE Overview. ROTATE is an open-ended learning framework for AHT. The core idea of ROTATE is to improve the AHT agent by iteratively generating diverse teammates with whom the AHT agent struggles to collaborate, yet not so adversarial that effective teamwork becomes impossible.

maximize different notions of diversity. One such diversity notion is *adversarial diversity* [42, 10], which seeks to generate a set of teams that cooperate well within teams, but not across teams. However, prior work [16, 47, 11] empirically demonstrates that adversarial diversity often leads to teammate policies that actively diminish returns when interacting with agents other than their identified teammate, a phenomenon sometimes called *self-sabotage*.

This paper addresses two issues that cause current methods to fail to learn policies that effectively collaborate with some teammates. First, two-stage AHT methods [40, 66, 41, 42] learn from interacting with teammates from a small fixed training set. Even when the training set is diverse, the AHT agent remains incapable of collaborating effectively with some teammates sampled from the vast space of possible strategies, specifically those with significantly different behavior from the policies in the training set [54, 41]. Second, other work designs a diverse training set of teammate policies by maximizing adversarial diversity [10, 42], which yields self-sabotaging teammates whose return-diminishing tendencies make it challenging for a randomly initialized RL-based AHT agent to learn to collaborate effectively [16, 47]. Despite addressing the first issue, some methods remain susceptible to the second issue by optimizing adversarial diversity [64].

We address the shortcomings of using a small, fixed training set by proposing an open-ended learning framework that continually generates new teammates with whom the AHT agent interacts to enhance its collaborative capabilities. We formulate our learning objective by observing that maximizing the expected returns of an AHT agent on a known set of teammates is equivalent to minimizing its expected *cooperative regret*: the utility gap between the best response to a given teammate, and the AHT agent’s performance with that teammate. While not knowing the teammates that will be encountered, we take inspiration from unsupervised environment design (UED) methods [57, 17, 24, 45] and train an AHT agent to minimize its regret against generated teammates that maximize the AHT agent’s cooperative regret. We propose a novel and practical objective that, unlike UED methods that optimize regret only at the initial state, also maximizes regret in states encountered later in an interaction. We build on these foundations to propose a practical algorithm, ROTATE (Fig. 1), which optimizes a cooperative regret-based minimax objective while maintaining a population of all teammates explored. We demonstrate that ROTATE significantly improves the robustness of AHT agents when faced with previously unseen teammates, compared to a range of baselines on two-player Level-Based Foraging and Overcooked tasks.

This paper makes three main contributions. First, it defines a novel problem formulation for AHT, enabling open-ended AHT training that continually generates new teammates. Second, it introduces a novel algorithm, ROTATE, that instantiates the proposed open-ended AHT framework. Third, it provides empirical evaluations demonstrating that ROTATE significantly improves return against unseen teammates compared to representative baselines from AHT and open-ended learning.

2 Related Work

Training AHT Agents. The training of ego agent policies that near-optimally collaborate with diverse previously unseen teammates has been explored in AHT [51]. Most AHT methods follow the two-stage design process, where the generation of a fixed training set of teammate policies is followed by AHT training. Given teammates from the training set, AHT methods [35] train an ego agent to model teammates [3] by first identifying their important characteristics (e.g., goals, beliefs, policies) based on their observed behavior, and then estimating the best-response policy to these teammates based on the inferred characteristics. Recent AHT methods [41, 40, 66, 56] use

neural networks trained using reinforcement learning [49, 36]. To further improve AHT training, several approaches learn a distribution for sampling teammate policies during training based on maximizing the worst-case returns [55] or regret [19, 12] of trained agents. While few, exceptions to the two-stage process include methods ~~that expect to encounter a continual stream of teammates at deployment~~ designed for continual AHT [38, 39, 64], ~~and methods that co-evolve populations of ego agents and teammates~~ [61, 64], self-play based methods, which do not explicitly optimize for diversity [62, 14], ~~and empirical game theoretic methods that optimize for cooperative diversity as a heuristic to induce generalization to unseen teammates~~ [29].

Teammate Generation for AHT & ZSC. Recent work removes the need to predefine teammate policy sets by generating diverse teammates during or before agent training. Other-Play [22] creates symmetry-equivalent teammates while training the agent policy, while E3T [62] mixes the agent’s current policy with a random policy to encourage diversity. FCP [52] trains teammates via repeated CMARL runs with different seeds, later improved by methods maximizing information-theoretic diversity objectives such as Jensen-Shannon divergence [34], mutual information [33], and entropy [60, 65]. More recent approaches [10, 43, 64] generate teammates that require distinct best-response strategies by maximizing adversarial diversity metrics, similar to ROTATE’s cooperative regret. Unlike ROTATE, these methods (i) maximize regret between generated teammates rather than with the trained agent, (ii) fix the teammate set prior to training, and (iii) evaluate regret only at the initial state. This last property leads to sabotaging teammates that harm cooperation in states unseen in self-play, motivating heuristic solutions in prior work [16, 47], and a systematic objective in ROTATE.

Open-Ended Learning (OEL). OEL [28, 53] studies algorithms that continually generate novel tasks to train generally capable agents [23, 4]. Many OEL approaches in RL take the form of unsupervised environment design (UED) [17], which improves generalization by designing or sampling new environments with varied initial states. Some methods directly train neural networks to propose environments that induce high regret in the agent [17], while others selectively sample curated tasks generated by procedural generators based on criteria such as expected return [57], TD-error [25], regret [24], or learnability [45]. In competitive MARL, OEL often produces new opponents through self-play [50, 30]. For AHT, MACOP [64] generates novel teammates via an adversarial diversity objective optimized with evolutionary methods and similar to the objectives studied by by Charakorn et al. [10] and Rahman et al. [42]. Thus, the objective can yield sabotaging teammates when applied only to the initial state. In contrast, ROTATE adopts a more systematic training objective that we demonstrate leads to performance gains.

3 Background

The interaction between agents in an AHT setting may be modeled as a decentralized Markov decision process (Dec-MDP) [5]. A Dec-MDP is characterized by a 7-tuple, $\langle N, S, \{\mathcal{A}^i\}_{i=1}^{|N|}, P, p_0, R, \gamma \rangle$, where N , S , and γ respectively denote the index set of agents within an interaction, the state space, and a discount rate, $0 \leq \gamma \leq 1$. Every interaction between agents begins from a state sampled from the initial state distribution, $s_0 \sim p_0(s)$. At timestep t , each agent, $i \in N$, jointly executes an action selected from its action space, $a_t^i \in \mathcal{A}^i$, based on the observed state, s_t , and its policy, $\pi^i(s_t^i)$. We assume that teammates choose their actions only based on the current state. Meanwhile, the AHT agent, also referred to as the *ego agent*, selects actions based on its state-action history, which is necessary to distinguish between different types of teammates effectively. Denoting the set of all probability distributions over a set X as $\Delta(X)$, the execution of the joint action, $a_t = (a_t^1, \dots, a_t^{|N|})$, results in agents observing a new state, s_{t+1} , sampled according to the environment transition function, $P : S \times \mathcal{A}^1 \times \dots \times \mathcal{A}^{|N|} \mapsto \Delta(S)$, and a common scalar reward, r_t , based on the reward function, $R : S \times \mathcal{A}^1 \times \dots \times \mathcal{A}^{|N|} \mapsto \mathbb{R}$.

4 Ad Hoc Teamwork Problem Formulation

AHT methods aim to train an adaptive policy that an ego agent can follow to achieve maximal return when collaborating with an unknown set of evaluation teammates. Formalizing the interaction between agents as a Dec-MDP, this section outlines the objective of AHT. While the most general AHT setting considers a possibly varying number of ego agents and teammates within an interaction [56, 41], this formalization addresses the more straightforward case where there is only a single ego agent within a team.

Let π^{ego} refer to the ego agent’s policy, and π^{-i} denote the $|N| - 1$ policies of the AHT agent’s teammates. We denote the returns of an ego agent that follows π^{ego} to collaborate with teammates controlled by π^{-i} , starting from state s , as:

$$V(s|\pi^{-i}, \pi^{\text{ego}}) = \mathbb{E}_{\substack{a_t^{\text{ego}} \sim \pi^{\text{ego}}, \\ a_t^{-i} \sim \pi^{-i}, P}} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \middle| s_0 = s \right]. \quad (1)$$

Let Π^{eval} denote the unknown set of joint teammate policies encountered during evaluation, which is assumed to only contain competent and non-adversarial policies, as defined in the seminal work of Stone et al. [51]. Let $\psi^{\text{eval}}(\Pi^{\text{eval}})$ denote the probability distribution over Π^{eval} defining how teammates are sampled during evaluation. An ego agent policy, π^{ego} , is evaluated by its ability to maximize the expected returns when collaborating with joint teammate policies sampled from $\psi^{\text{eval}}(\Pi^{\text{eval}})$, which is formalized as:

$$\max_{\pi^{\text{ego}}} V(\psi^{\text{eval}}, \Pi^{\text{eval}}, \pi^{\text{ego}}) = \max_{\pi^{\text{ego}}} \mathbb{E}_{\pi^{-i} \sim \psi^{\text{eval}}(\Pi^{\text{eval}}), s_0 \sim p_0} [V(s_0|\pi^{-i}, \pi^{\text{ego}})]. \quad (2)$$

An optimal π^{ego} that maximizes Eq. 2 closely approximates the *best response policy* performance when collaborating with $\pi^{-i} \in \Pi^{\text{eval}}$. Given a teammate policy π^{-i} , $\text{BR}(\pi^{-i})$ is a best response policy to π^{-i} if and only the team policy formed by π^{-i} and $\text{BR}(\pi^{-i})$ achieves maximal return:

$$\text{BR}(\pi^{-i}) \in \arg \max_{\pi} \mathbb{E}_{s \sim p_0} [V(s|\pi, \pi^{-i})]. \quad (3)$$

In some cases, AHT algorithms can estimate this optimal policy by using Π^{eval} to train an ego agent policy that maximizes $V(\psi^{\text{eval}}, \Pi^{\text{eval}}, \pi^{\text{ego}})$ when Π^{eval} is known.¹ However, most AHT methods address the more challenging case where Π^{eval} is unknown, which is the setting that this paper adopts as well. While our methods assume no knowledge of Π^{eval} during training, we follow standard practice [40, 41, 66, 56] by manually designing a diverse Π^{eval} for evaluation purposes, as we later describe in Section 7.

When Π^{eval} is unknown, AHT algorithms [35] learn by interacting with policies from the training set, Π^{train} , which are learned or manually designed by leveraging an expert’s domain knowledge about the characteristics of Π^{eval} . After forming the set of training teammates, current AHT algorithms use RL to discover an ego agent policy based on interactions with joint policies sampled from Π^{train} . While the precise training objective varies with the AHT algorithm, methods commonly estimate the ego agent policy maximizing the expected return during interactions with joint policies sampled uniformly from Π^{train} , which we describe below:

$$\pi^{*,\text{ego}}(\Pi^{\text{train}}) = \arg \max_{\pi^{\text{ego}}} \mathbb{E}_{\pi^{-i} \sim \mathcal{U}(\Pi^{\text{train}}), s_0 \sim p_0} [V(s_0|\pi^{-i}, \pi^{\text{ego}})]. \quad (4)$$

Naturally, even $\pi^{*,\text{ego}}(\Pi^{\text{train}})$ may be suboptimal with respect to Π^{eval} and ψ^{eval} , due to the potential distribution shift caused by differences between the training and evaluation objectives.

5 An Open-Ended Learning Perspective on Reformulating Ad Hoc Teamwork as an Open-Ended Learning Problem

In this section, we outline the general components of our show how the idealized ad hoc teamwork objective—training ego agents to collaborate well with unknown teammates (Eq. 2, Section 4)—can be operationalized as a cooperative regret-driven, open-ended framework to train ego agents that are performant at collaborating with holdout teammate policies, despite not knowing Π^{eval} and ψ^{eval} during training. We first argue for minimizing worst-case cooperative regret towards training ego agent policies that maximize learning procedure. In particular, we show that for a fixed set of teammates Π^{eval} and sampling distribution ψ^{eval} over Π^{eval} , maximizing the return of the ego agent is equivalent to minimizing its cooperative regret. In absence of knowledge about Π^{eval} and ψ^{eval} , we argue that minimizing the worst-case cooperative regret of the ego agent with respect to regret maximizing teammates leads to ego agents that cooperate well with any unknown teammate. Based on this, we propose a novel *minimax regret* objective (Eq. 2 when Π^{eval} is unknown. We then

¹In the context of reinforcement-learning-based AHT algorithms, “known” means that an AHT algorithm has unlimited sampling access to the teammate policies.

finish the section by introducing two necessary procedures 7). Finally, we present an algorithmic framework for optimizing the minimax regret objective in an iterative process to minimize worst-case regret fashion (Alg. 1).

We define the *cooperative regret* of an ego agent policy π^{ego} when interacting with some joint teammate policy π^{-i} from a starting state s as:

$$\text{CR}(\pi^{\text{ego}}, \pi^{-i}, s) = V(s | \pi^{-i}, \text{BR}(\pi^{-i})) - V(s | \pi^{-i}, \pi^{\text{ego}}). \quad (5)$$

Any optimal AHT policy that maximizes Eq. 2 must also minimize the expected regret over joint teammate policies sampled based on $\psi^{\text{eval}}(\Pi^{\text{eval}})$, which we formally express as:

$$\text{CR}(\psi^{\text{eval}}, \Pi^{\text{eval}}, \pi^{\text{ego}}) = \mathbb{E}_{\pi^{-i} \sim \psi^{\text{eval}}(\Pi^{\text{eval}}), s_0 \sim p_0} [\text{CR}(\pi^{\text{ego}}, \pi^{-i}, s_0)]. \quad (6)$$

This property is a consequence of $V(s | \pi^{-i}, \text{BR}(\pi^{-i}))$ being independent of π^{ego} for any π^{-i} and s , leaving maximizing expected regret equivalent to minimizing the negative expected returns when collaborating with joint teammate policies sampled from $\psi^{\text{eval}}(\Pi^{\text{eval}})$.

Without knowing Π^{eval} to optimize $\text{CR}(\psi^{\text{eval}}, \Pi^{\text{eval}}, \pi^{\text{ego}})$, we instead take inspiration from approaches in UED [57, 17], and propose optimizing π^{ego} to minimize the *worst-case regret* that could be induced by any teammate policy π^{-i} :

$$\min_{\pi^{\text{ego}}} \max_{\pi^{-i} \in \Pi^{-i}} \mathbb{E}_{s_0 \sim p_0} [\text{CR}(\pi^{\text{ego}}, \pi^{-i}, s_0)], \quad (7)$$

where we re-emphasize that Π^{-i} denotes the set of all competent and non-adversarial [51] joint teammate policies. Limiting the considered joint policies is important, as teams that consistently perform poorly against any π^{ego} are unlikely to be encountered in coordination scenarios and may introduce unnecessary learning challenges for RL-based AHT learning algorithms.

Finding π^{ego} that achieves zero worst-case regret is equivalent to finding an ego agent that achieves the best-response return with any joint teammate policy π^{-i} . If such a π^{ego} exists, then this AHT agent would maximize Eq. 2 for any ψ^{eval} and Π^{eval} . However, its existence is not guaranteed [31]. In practice, we are content with *minimizing* the worst-case regret. While minimizing worst-case regret still applies to AHT problems with more than one teammate, we limit our method for optimizing Eq. 7 and our experiments to two-player, fully observable AHT games.

Algorithm 1 Open-Ended Ad Hoc Teamwork Framework

Require:

Environment, Env.
Total of training iterations, T^{iter}
Initial ego agent policy parameters, θ^{ego}

| | |
|---|---|
| <pre> 1: $\mathbf{B}_{\pi} \leftarrow \langle \rangle$ 2: for $j = 1, 2, \dots, T^{\text{iter}}$ do 3: $\mathbf{B}_{\pi}^{\text{new}} \leftarrow \text{TeammateGenerator}(\text{Env}, \theta^{\text{ego}}, \mathbf{B}_{\pi})$ 4: $\theta^{\text{ego}} \leftarrow \text{EgoUpdate}(\text{Env}, \theta^{\text{ego}}, \mathbf{B}_{\pi}^{\text{new}})$ 5: $\mathbf{B}_{\pi} \leftarrow \mathbf{B}_{\pi}^{\text{new}}$ 6: end for 7: Return θ^{ego} </pre> | <p>▷ Init teammate policy parameter buffer.</p> <p>▷ Train teammates to maximize regret.</p> <p>▷ Train ego agent to minimize regret.</p> |
|---|---|

Algorithm 1 in the Appendix outlines a framework to train a π^{ego} that minimizes outlines our general framework for training an ego agent to minimize the worst-case regret. The algorithm cooperative regret induced by any teammate $\pi^{-i} \in \Pi^{-i}$. Algorithm 1 resembles coordinate ascent algorithms [18], which alternate between optimizing for π^{-i} and π^{ego} for T^{iter} iterations, while assuming the other is fixed. We call a phase where we fix π^{ego} and update π^{-i} to maximize the ego agent's regret, the *teammate generation phase*. Meanwhile, assuming that π^{-i} is fixed, the *ego agent update phase* updates π^{ego} to minimize regret. This optimization algorithm is

Our practical algorithm, ROTATE, instantiates Algorithm 1 by specifying the **TeammateGenerator** and **EgoUpdate** procedures, and is described in Section 6. It is an open-ended training method that continually generates novel teammate policies whose interaction with the ego agent provides the learning experience to improve π^{ego} . Next, we detail the learning process during these two

phases learning procedure according to the definition proposed by Hughes et al. [23], because it continually generates novel yet learnable artifacts (i.e., teammates) for an observer (i.e., ego agent). A discussion of how ROTATE satisfies the definition of Hughes et al. [23] is provided in App. C.1.

6 Regret-driven Open-ended Training for Ad Hoc Teamwork Practical Algorithm: ROTATE

This section presents our regret-driven practical algorithm for optimizing the minimax regret objective proposed in Section 5, open-ended AHT algorithm, ROTATE. We first describe the teammate generation procedure in Section 6.1, focusing on motivating the objective used to generate teammate policies. Next, we describe the ego agent update method in Section 6.2. App. A provides ROTATE’s the ROTATE pseudocode and a more detailed implementation description discussion of the losses and exact update procedure.

6.1 ROTATE Teammate Generator

Given a fixed π^{ego} , ROTATE’s teammate generator seeks to discover a teammate policy that maximizes the cooperative regret of π^{ego} . Maximizing cooperative regret requires estimating the teammate policy, π^{-i} , and its best response policy, $\text{BR}(\pi^{-i})$. In the following, we abbreviate $\text{BR}(\pi^{-i})$ to BR for brevity. ROTATE’s teammate generator estimates both policies using the Proximal Policy Optimization (PPO) algorithm [49].

The *per-trajectory regret* of π^{ego} (i.e., the inner objective of Eq. 7) is the regret from trajectories starting from the initial state distribution:

$$\max_{\pi^{-i}} \mathbb{E}_{s_0 \sim p_0} [\text{CR}(\pi^{\text{ego}}, \pi^{-i}, s_0)]. \quad (8)$$

Eq. 8 resembles the objectives used in past UED [57, 17] and the teammate generation literature [43, 10] to generate tasks or teammate policies. Recent work demonstrates that maximizing per-trajectory regret is prone to yielding self-sabotaging teammates [16, 47]. Maximizing the cooperative regret only from $s_0 \sim p_0$ implicitly encourages $\text{BR}(\pi^{-i})$ to select actions leading to future states that are distinguishable from those encountered during the interaction between π^{-i} and π^{ego} . When encountering future states from interactions with π^{ego} , π^{-i} ends up choosing actions that sabotage cooperation by minimizing the teammate’s returns against π^{-i} . Thus, training π^{ego} to minimize regret (i.e., by maximizing the expected returns) when collaborating with π^{-i} using RL becomes challenging because π^{-i} actively chooses actions that undermine collaboration.

We mitigate the emergence of self-sabotage by training π^{-i} to maximize two objectives across states sampled from different state visitation distributions. These state visitation distributions result from: (i) teammate and BR interactions (self-play, SP), (ii) teammate and ego agent interactions (cross-play, XP), and (iii) cross-play continued by self-play interactions (SXP), where the teammate is first interacting with the ego agent, but switches at a random timestep t to interacting with its BR. Let $d(\pi^1, \pi^2; p)$ denote the state visitation distribution when π^1 and π^2 interact based on a starting state distribution p . We use the following shorthand to denote the SP, XP, and SXP state visitation distributions:

$$p_{\text{SP}} := d(\pi^{-i}, \text{BR}(\pi^{-i}); p_0), \quad p_{\text{XP}} := d(\pi^{-i}, \pi^{\text{ego}}; p_0), \quad p_{\text{SXP}} := d(\pi^{-i}, \text{BR}(\pi^{-i}); p_{\text{XP}}). \quad (9)$$

Based on these distributions, we define the following *per-state regret* objective for training π^{-i} :

$$\max_{\pi^{-i}} (\mathbb{E}_{s \sim 0.5 p_{\text{XP}} + 0.5 p_{\text{SP}}} [\text{CR}(\pi^{\text{ego}}, \pi^{-i}, s)] + \mathbb{E}_{s \sim p_{\text{SXP}}} [V(s | \pi^{-i}, \text{BR}(\pi^{-i}))]). \quad (10)$$

The difference between the per-trajectory and per-step regret objectives is visualized in Figure 2. Both terms in the per-state regret objective discourage adversarial behavior from π^{-i} . The first term in Expr. 10 corresponds to the ego agent’s regret starting from both SP and XP states. Estimating regret from XP and SP requires collecting SXP data as well as an analogous type of data called XSP (SP continued by XP interactions), as detailed in App. A.2. In general, optimizing the ego agent’s regret encourages discovering π^{-i} for which the ego agent policy has a high room for improvement. Optimizing regret starting from XP states requires π^{-i} to be able to coordinate with its BR starting from any state encountered during interactions with the ego agent, thus preventing π^{-i} from irrecoverably sabotaging an interaction. On the other hand, optimizing regret from SP states

requires π^{-i} to be able to decrease the return of the ego agent starting from any state encountered during interactions between the teammate and the BR, thus disincentivizing the emergence of unconditional cooperation signals. Finally, we find that training π^{-i} to collaborate well with its BR even during SXP interactions helps ensure that π^{-i} is a good-faith collaborator with at least one partner.

While obtaining states from p_{SP} and p_{XP} is straightforward, states from p_{SXP} and p_{XSP} are collected using either environment re-setting or policy switching. Using SXP as an example, if an environment supports re-setting to any arbitrary state, then states encountered during XP interaction can be stored and used as the initial state for SP interactions. Otherwise, we may sample a random timestep t , run XP interaction until timestep t , and then switch to SP interaction [47]. Only data gathered after timestep t should be used to compute objectives based on p_{SXP} .

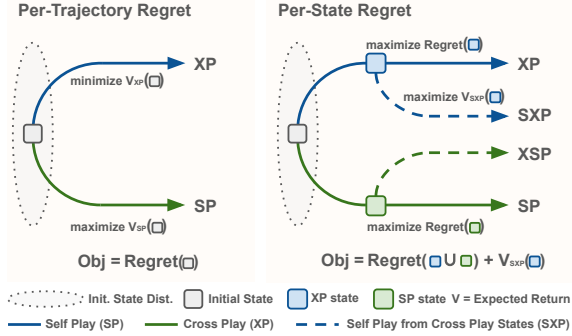


Figure 2: Teammate policy optimization objectives: per-trajectory regret vs per-state regret.

6.2 ROTATE Ego Agent Update

At each iteration, ROTATE creates a teammate that attempts to discover cooperative weaknesses of the previous iteration’s ego agent, by maximizing its per-state regret. To allow the ROTATE ego agent to improve its robustness over time and reduce the possibility that it forgets how to cope with previously generated teammates, the ROTATE ego agent maintains a *population buffer* of generated teammates. During the ego agent update phase of each iteration, the ROTATE ego agent is trained using PPO [49] against teammates sampled uniformly from the population buffer. We find experimentally that for the ego agent to learn effectively against the nonstationary population buffer, it is important to define a lower entropy coefficient and learning rate than when training the teammate and BR agents (typically in the range of 1×10^{-4} for the entropy coefficient and 1×10^{-5} for the learning rate).

7 Experimental Results

This section presents the empirical evaluation of ROTATE compared to baseline methods, [across as well as several ablations](#). [The experiments consider](#) one illustrative matrix game and six benchmark tasks. Supplemental results, implementation details, and code link are provided in the Appendix. The main research questions are:

- **RQ1:** Does ROTATE better generalize to unseen teammates, compared to baseline methods from the AHT and UED literature? (Yes)
- **RQ2:** Does per-state regret mitigate sabotage and improve generalization to unseen teammates compared to per-trajectory regret? (Yes)
- **RQ3:** [Does the SXP return term of the ROTATE teammate generation objective improve learning and generalization?](#) (Yes)
- **RQ4:** [Is the population buffer necessary for ROTATE to learn well?](#) (Yes)

7.1 Experimental Setup

This section describes the experimental setting, including tasks, baselines, construction of the evaluation set, and the evaluation metric.

Tasks ROTATE is evaluated on a didactic matrix game and six benchmark tasks. For clarity, the matrix game is described with the corresponding results. The benchmark tasks are Level-Based Foraging (LBF) [2], and the five classic layouts from the Overcooked suite [9]: Cramped Room (CR), Asymmetric Advantages (AA), Counter Circuit (CC), Coordination Ring (CoR), and Forced Coordination (FC). All six tasks are cooperative, permit a variety of possible conventions, and are commonly used within the AHT literature [2, 13, 40]. In LBF, two agents must navigate to apples that are randomly placed within a gridworld, and cooperate to pick up the apples. In all Overcooked

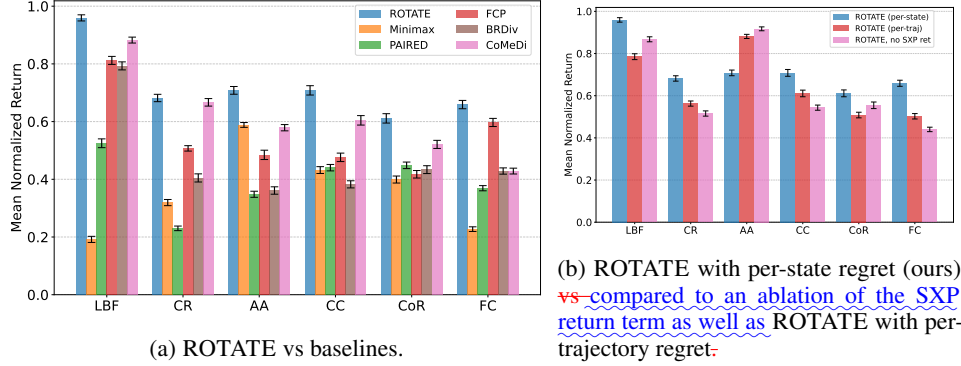


Figure 3: (Left) ROTATE outperforms all baseline methods across all tasks in evaluation return. (Right) ROTATE with per-state regret (ours) outperforms ROTATE with per-trajectory regret in 5/6 tasks. 95% bootstrapped CI’s are shown, computed across all evaluation teammates and trials.

tasks, two agents collaborate in varying gridworld kitchen layouts to prepare dishes. All experiments were implemented with JAX [8].

Baselines As our method is most closely related to methods from UED and teammate generation, we compare against two UED methods adapted for AHT (PAIRED [17], Minimax Return [37, 55]) and three teammate generation methods (Fictitious Co-Play [52], BRDiv [42], CoMeDi [47]). While curator-based methods such as PLR [24, 25] are prevalent in UED, we do not compare against them as they are orthogonal to ROTATE [19, 55, 12]. Similarly, we do not compare against AHT algorithms for ego learning [3]. Each baseline is described in detail in App. B. For fair comparison, all open-ended and UED methods were trained for a similar number of environment interactions, or until best performance on the evaluation set. All teammate generation approaches were ran using a similar number of environment interactions as their original implementations, as scaling them up to use a similar number of steps as the open-ended approaches proved challenging (see discussion in App. B). All results are reported with three trials.

Construction of Π^{eval} We wish to evaluate all methods on as diverse a set of evaluation teammates as practically feasible, while ensuring that each teammate acts in “good faith”. To achieve this goal, for each task, we construct 9 to 13 evaluation teammates using three methods: IPPO with varied seeds and reward shaping, BRDiv, and manually programmed heuristic agents. Full descriptions of the teammate generation procedure and all teammates in Π^{eval} are provided in App. G.

Evaluation Metric Ego agent policies are evaluated with each teammate in Π^{eval} for 64 evaluation episodes, where the return is computed for each episode, and normalized using a lower return bound of zero and an estimated best response return as the upper bound for each teammate. Performance of a method on Π^{eval} is reported as the normalized mean return with bootstrapped 95% confidence intervals, computed via the `rliable` library [1]. Our normalized return metric is similar to the BRProx metric recommended by Wang et al. [58]. Details about the normalization procedure and specific bounds for each teammate are reported in the App. G.

7.2 Results

This section addresses the ~~three~~ research questions introduced at the beginning of Section 7. Supplemental analysis considering alternative regret-based objectives, independent utility of the ROTATE population, performance breakdowns by evaluation teammate, learning curves for all variants of ROTATE, and a human proxy evaluation on Overcooked, are provided in App. D.

RQ1: Does ROTATE better generalize to unseen teammates compared to baselines? (Yes) We evaluate ROTATE’s generalization capabilities by comparing its performance against baselines on Π^{eval} . Fig. 3a compares the normalized mean returns for ROTATE and baseline methods across the six tasks. The results show that ROTATE significantly outperforms all baselines on 5/6 tasks.

Among the baseline methods, the next best performing baselines are CoMeDi and FCP. We attribute CoMeDi’s strong performance to the resemblance of its mixed-play objective to our per-state regret objective, which we discuss in App. C.3. FCP’s performance may be attributed to the large number of partners that FCP was trained with (approximately 100 teammates per task). We found that FCP

tends to perform especially well with the IPPO policies in Π^{eval} , likely because the IPPO evaluation teammates are in-distribution for the distribution of teammates constructed by FCP.

Minimax Return performs surprisingly well in AA, which may be attributed to AA’s particular characteristics. In AA, agents operate in separated kitchen halves, possessing all necessary resources for individual task completion, with pots on the dividing counter being the only shared resource. A team where both agents act fully independently may achieve high returns—albeit coordination leads to still higher returns.² Visualizing the trained policies reveals that the adversarial teammate trained by Minimax Return cannot drive the ego agent’s return to zero, and does not prevent the ego agent from learning how to perform the task independently. However, on LBF and FC, where coordination is crucial to obtain positive returns, Minimax Return is the worst-performing baseline.

BRDiv and PAIRED exhibit comparatively poor performance, which may be partially attributed to their teammate generation objectives that resemble per-trajectory regret. As we find for **RQ2**, per-state regret outperforms per-trajectory regret within the ROTATE framework. Furthermore, PAIRED’s update structure involves a lockstep training process for the teammate generator, best response, and ego agent. This synchronized training may hinder the natural emergence of robust conventions that are crucial for effective AHT.

| | H | T | S |
|---|----|----|----|
| H | 1 | 0 | -1 |
| T | 0 | 1 | -1 |
| S | -1 | -1 | -1 |

Table 1: Payoff matrix for the sabotage game.

RQ2a: Does per-state regret mitigate sabotage compared to per-trajectory regret? (Yes) We design a simple *sabotage game* to investigate the whether the per-state regret objective leads to teammate policies that sabotage less often compared to the per-trajectory regret objective. The sabotage game is a fully cooperative, iterated matrix game with payoff matrix shown in Table 1. Each agent observes a game state that consists of the complete history of joint actions. Both agents have three actions: H, T, and S(abotage). The first two actions lead to two possible cooperative outcomes, while the last action leads to a reward of -1 if selected by either agent and immediate episode termination. Thus, the last action corresponds to sabotaging the team’s payoffs. By default, the game lasts for five timesteps.

We train ROTATE with both per-state and per-trajectory regret. To measure the extent to which the learned teammate policies engage in sabotage, we enumerate the 341 non-terminal states in the game and measure the probability of the sabotage action at each state for the last generated teammate policy. Fig. 4 shows that ROTATE with per-state regret has a near-zero probability of taking the sabotage action at all non-terminal states, while the per-trajectory regret objective leads to over a third of states that have $P(S)$ near 1.0.

RQ2b: Does per-state regret lead to improved generalization compared to per-trajectory regret? (Yes) RQ2a demonstrated that ROTATE with per-state regret (ours) leads to teammate policies that sabotage less often in an illustrative matrix game, compared to ROTATE with per-trajectory regret. Here, we investigate whether the this translates to improved generalization against the unseen evaluation teammates. All configurations other than the teammate’s policy objective are kept identical, including the data used to train the teammate value functions. Fig. 3b shows that ROTATE with per-state regret outperforms ROTATE with per-trajectory regret on all tasks except AA, confirming the superiority of per-state regret. As discussed in **RQ1**, we observe that AA is a layout where an ego agent is less susceptible to sabotage, due to the separated kitchen layout. App. ??-C.4 presents additional experiments testing ROTATE with CoMeDi-style mixed-play rollouts, and alternative methods to compute per-state regret—ultimately finding that ROTATE outperforms all variations.

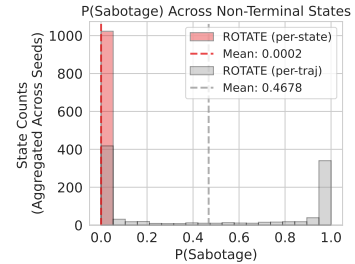


Figure 4: Probability of the sabotage action at all states in the sabotage game for ROTATE teammates trained with per-state regret (ours) vs per-trajectory regret. Results are aggregated across three trials.

²Optimal behavior in AA still requires effective coordination due to layout asymmetry. In the “left” kitchen, the delivery zone is adjacent to the pots while the onions are farther, while in the “right” kitchen, the opposite is true. Thus, an optimal team consists of the “left” agent delivering finished soup, and the “right” agent placing onions in the pots—and indeed, we observe that teams of IPPO agents converge to this behavior.

RQ3: Does the SXP return term of the ROTATE teammate generation objective improve learning and generalization? (Yes) The ROTATE objective (Eq. 10) includes two terms: the per-state regret term and a SXP return term, introduced to ensure that the teammate collaborates well with its BR, even during SXP interactions. This helps ensure that the teammate is a good-faith collaborator. In Fig. 3b, the heldout return of an ablation of ROTATE where the SXP return term is removed is shown. We find that the ablated version performs significantly worse than the non-ablated across 5/6 of the tasks, confirming that the auxiliary return term plays an important role in the teammate generation objective.

RQ3RQ4: Is the population buffer necessary for ROTATE to learn well? (Yes) We hypothesize that collecting all previously generated teammates in a population buffer helps the ROTATE agent improve in robustness against all previously discovered conventions. On the other hand, if there is no population buffer, then it becomes possible for the ROTATE ego agent to forget how to collaborate with teammate seen at earlier iterations of open-ended learning [27], which creates the possibility that the ego agent and teammate generator oscillates between conventions. As shown in Fig. 6a, ROTATE without the population buffer attains lower evaluation returns than the full ROTATE method on all tasks except for AA, thus supporting the hypothesis that the population buffer improves ego agent learning. As discussed in **RQ1**, AA is a unique layout where agents can complete the task independently, even in the presence of an adversarial partner. As a corollary, there are few meaningful cooperative conventions that can be discovered, and no scenarios where convention mismatch leads to zero return (unlike LBF and FC).

8 Discussion and Conclusion

This paper reformulates AHT as an open-ended learning problem and introduces ROTATE, a regret-driven algorithm. ROTATE iteratively alternates between improving an AHT agent and generating challenging yet cooperative teammates by optimizing a per-state regret objective designed to discover teammates that exploit cooperative vulnerabilities while mitigating self-sabotage. Experiments on an illustrative matrix game demonstrate that the per-state regret objective mitigates self-sabotage. Extensive evaluations across six benchmark tasks demonstrate that ROTATE significantly enhances the generalization capabilities of AHT agents when faced with previously unseen teammates, outperforming baselines from both AHT and UED.

The current work has several limitations. First, while this paper provides intuitive justification and strong empirical evidence for the efficacy of the per-state regret objective, an exciting line of follow-up work is to formally define the concept of self-sabotage and theoretically analyze the properties of the proposed regret objectives. Second, the paper only validates ROTATE on two-agent, fully observable, and fully cooperative scenarios, which leaves the question of whether it scales to more complex scenarios for future work. Finally, this work has focused on the teammate generation phase of open-ended AHT. Future work might explore ego agent training methods that better handle the nonstationarity induced by open-ended teammate generation.

References

- [1] R. Agarwal, M. Schwarzer, P. S. Castro, A. C. Courville, and M. Bellemare. Deep Reinforcement Learning at the Edge of the Statistical Precipice. In *Advances in Neural Information Processing Systems*, volume 34, pages 29304–29320. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/hash/f514cec81cb148559cf475e7426eed5e-Abstract.html.
- [2] S. V. Albrecht and S. Ramamoorthy. A game-theoretic model and best-response learning method for ad hoc coordination in multiagent systems. In *Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems*, AAMAS '13, pages 1155–1156, Richland, SC, May 2013. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 978-1-4503-1993-5.
- [3] S. V. Albrecht and P. Stone. Autonomous agents modelling other agents: A comprehensive survey and open problems. *Artificial Intelligence*, 258:66–95, 2018. ISSN 0004-3702. doi: <https://doi.org/10.1016/j.artint.2018.01.002>. URL <https://www.sciencedirect.com/science/article/pii/S0004370218300249>.

- [4] B. Baker, I. Kanitscheider, T. Markov, Y. Wu, G. Powell, B. McGrew, and I. Mordatch. Emergent Tool Use From Multi-Agent Autocurricula. In *International Conference on Learning Representations*, Sept. 2019. URL <https://openreview.net/forum?id=SkxpxJBKwS>.
- [5] D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein. The complexity of decentralized control of markov decision processes. *Mathematics of operations research*, 27(4):819–840, 2002.
- [6] C. Bonnet, D. Luo, D. J. Byrne, S. Surana, S. Abramowitz, P. Duckworth, V. Coyette, L. I. Midgley, E. Tegegn, T. Kalloniatis, O. Mahjoub, M. Macfarlane, A. P. Smit, N. Grinsztajn, R. Boige, C. N. Waters, M. A. A. Mimouni, U. A. M. Sob, R. J. d. Kock, S. Singh, D. Furelos-Blanco, V. Le, A. Pretorius, and A. Laterre. Jumanji: a Diverse Suite of Scalable Reinforcement Learning Environments in JAX. In *The Twelfth International Conference on Learning Representations*, Oct. 2023. URL <https://openreview.net/forum?id=C4CxQmp9wc>.
- [7] M. Bowling and P. McCracken. Coordination and adaptation in impromptu teams. In *AAAI*, volume 5, pages 53–58, 2005.
- [8] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/jax-ml/jax>.
- [9] M. Carroll, R. Shah, M. K. Ho, T. Griffiths, S. Seshia, P. Abbeel, and A. Dragan. On the Utility of Learning about Humans for Human-AI Coordination. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/hash/f5b1b89d98b7286673128a5fb112cb9a-Abstract.html.
- [10] R. Charakorn, P. Manoonpong, and N. Dilokthanakul. Generating diverse cooperative agents by learning incompatible policies. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=UkU05G0H7_6.
- [11] R. Charakorn, P. Manoonpong, and N. Dilokthanakul. Diversity is not all you need: Training a robust cooperative agent needs specialist partners. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 56401–56423. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/66b35d2e8d524706f39cc21f5337b002-Paper-Conference.pdf.
- [12] P. Chaudhary, Y. Liang, D. Chen, S. S. Du, and N. Jaques. Improving human-ai coordination through adversarial training and generative models, 2025. URL <https://arxiv.org/abs/2504.15457>.
- [13] F. Christianos, L. Schäfer, and S. V. Albrecht. Shared experience actor-critic for multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [14] D. Cornelisse and E. Vinitsky. Human-compatible driving partners through data-regularized self-play reinforcement learning, June 2024. URL <http://arxiv.org/abs/2403.19648>.
- [15] B. Cui, H. Hu, L. Pineda, and J. Foerster. K-level reasoning for zero-shot coordination in hanabi. *Advances in Neural Information Processing Systems*, 34:8215–8228, 2021.
- [16] B. Cui, A. Lupu, S. Sokota, H. Hu, D. J. Wu, and J. N. Foerster. Adversarial diversity in hanabi. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=uLE3WF3-H_5.
- [17] M. Dennis, N. Jaques, E. Vinitsky, A. Bayen, S. Russell, A. Critch, and S. Levine. Emergent Complexity and Zero-shot Transfer via Unsupervised Environment Design. In *Advances in Neural Information Processing Systems*, volume 33, pages 13049–13061. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/985e9a46e10005356bbaf194249f6856-Abstract.html>.

- [18] D. d’Esopo. A convex programming procedure. *Naval Research Logistics Quarterly*, 6(1): 33–42, 1959.
- [19] H. Erlebach and J. Cook. RACCOON: Regret-based adaptive curricula for cooperation. In *Coordination and Cooperation for Multi-Agent Reinforcement Learning Methods Workshop*, 2024. URL <https://openreview.net/forum?id=jAH5JNY3Qd>.
- [20] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [21] P. W. Glynn. Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 33(10):75–84, 1990.
- [22] H. Hu, A. Lerer, A. Peysakhovich, and J. Foerster. “other-play” for zero-shot coordination. In *International Conference on Machine Learning*, pages 4399–4410. PMLR, 2020.
- [23] E. Hughes, M. D. Dennis, J. Parker-Holder, F. Behbahani, A. Mavalankar, Y. Shi, T. Schaul, and T. Rocktäschel. Position: Open-Endedness is Essential for Artificial Superhuman Intelligence. In *Proceedings of the 41st International Conference on Machine Learning*, pages 20597–20616. PMLR, July 2024. URL <https://proceedings.mlr.press/v235/hughes24a.html>.
- [24] M. Jiang, M. Dennis, J. Parker-Holder, J. Foerster, E. Grefenstette, and T. Rocktäschel. Replay-guided adversarial environment design. *Advances in Neural Information Processing Systems*, 34:1884–1897, 2021.
- [25] M. Jiang, E. Grefenstette, and T. Rocktäschel. Prioritized level replay. In *International Conference on Machine Learning*, pages 4940–4950. PMLR, 2021.
- [26] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In *3rd International Conference for Learning Representations*, San Diego, CA, 2015. doi: 10.48550/arXiv.1412.6980. URL <http://arxiv.org/abs/1412.6980>. arXiv:1412.6980 [cs].
- [27] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, Mar. 2017. doi: 10.1073/pnas.1611835114. URL <https://www.pnas.org/doi/10.1073/pnas.1611835114>. Publisher: Proceedings of the National Academy of Sciences.
- [28] W. Langdon. Pfeiffer—a distributed open-ended evolutionary system. In *AISB*, volume 5, pages 7–13. Citeseer, 2005.
- [29] Y. Li, S. Zhang, J. Sun, Y. Du, Y. Wen, X. Wang, and W. Pan. Cooperative Open-ended Learning Framework for Zero-Shot Coordination. In *Proceedings of the 40th International Conference on Machine Learning*, pages 20470–20484. PMLR, July 2023. URL <https://proceedings.mlr.press/v202/li23au.html>. ISSN: 2640-3498.
- [30] F. Lin, S. Huang, T. Pearce, W. Chen, and W.-W. Tu. TiZero: Mastering Multi-Agent Football with Curriculum Learning and Self-Play. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, AAMAS ’23, pages 67–76, Richland, SC, May 2023. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 978-1-4503-9432-1.
- [31] R. Loftin and F. A. Oliehoek. On the impossibility of learning to cooperate with adaptive partner strategies in repeated games. In *International Conference on Machine Learning*, pages 14197–14209. PMLR, 2022.
- [32] C. Lu, Y. Schroecker, A. Gu, E. Parisotto, J. N. Foerster, S. Singh, and F. Behbahani. Structured State Space Models for In-Context Reinforcement Learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, Nov. 2023. URL <https://openreview.net/forum?id=4W9FVg1j6I¬eId=38Anv4M4TW>.

- [33] K. Lucas and R. E. Allen. Any-play: An intrinsic augmentation for zero-shot coordination. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, pages 853–861, 2022.
- [34] A. Lupu, B. Cui, H. Hu, and J. Foerster. Trajectory diversity for zero-shot coordination. In *International conference on machine learning*, pages 7204–7213. PMLR, 2021.
- [35] R. Mirsky, I. Carlucho, A. Rahman, E. Fosong, W. Macke, M. Sridharan, P. Stone, and S. V. Albrecht. A survey of ad hoc teamwork research. In *European conference on multi-agent systems*, pages 275–293. Springer, 2022.
- [36] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PmLR, 2016.
- [37] J. Morimoto and K. Doya. Robust Reinforcement Learning. *Neural Comput.*, 17(2):335–359, Feb. 2005. ISSN 0899-7667. doi: 10.1162/0899766053011528. URL <https://doi.org/10.1162/0899766053011528>.
- [38] H. Nekoei, A. Badrinaaraayanan, A. Courville, and S. Chandar. Continuous Coordination As a Realistic Scenario for Lifelong Learning. In *Proceedings of the 38th International Conference on Machine Learning*, pages 8016–8024. PMLR, July 2021. URL <https://proceedings.mlr.press/v139/nekoei21a.html>. ISSN: 2640-3498.
- [39] H. Nekoei, X. Zhao, J. Rajendran, M. Liu, and S. Chandar. Towards Few-shot Coordination: Revisiting Ad-hoc Teamplay Challenge In the Game of Hanabi. In *Proceedings of The 2nd Conference on Lifelong Learning Agents*, pages 861–877. PMLR, Nov. 2023. URL <https://proceedings.mlr.press/v232/nekoei23b.html>. ISSN: 2640-3498.
- [40] G. Papoudakis, F. Christianos, and S. V. Albrecht. Agent modelling under partial observability for deep reinforcement learning. In *Advances in Neural Information Processing Systems*, 2021.
- [41] A. Rahman, N. Höpner, F. Christianos, and S. V. Albrecht. Towards Open Ad Hoc Teamwork Using Graph-based Policy Learning. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139. PMLR, June 2021.
- [42] A. Rahman, E. Fosong, I. Carlucho, and S. V. Albrecht. Generating teammates for training robust ad hoc teamwork agents via best-response diversity. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=l5BzfQhR0l>.
- [43] M. Rahman, J. Cui, and P. Stone. Minimum coverage sets for training robust ad hoc teamwork agents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17523–17530, 2024.
- [44] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21:1–51, 2020.
- [45] A. Rutherford, M. Beukman, T. Willi, B. Lacerda, N. Hawes, and J. N. Foerster. No regrets: Investigating and improving regret approximations for curriculum discovery. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=iEeiZlTbts>.
- [46] A. Rutherford, B. Ellis, M. Gallici, J. Cook, A. Lupu, G. Ingvarsson, T. Willi, R. Hammond, A. Khan, C. S. de Witt, A. Souly, S. Bandyopadhyay, M. Samvelyan, M. Jiang, R. Lange, S. Whiteson, B. Lacerda, N. Hawes, T. Rocktäschel, C. Lu, and J. Foerster. JaxMARL: Multi-Agent RL Environments and Algorithms in JAX. In *Advances in Neural Information Processing Systems*, volume 37, pages 50925–50951, Dec. 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/5aee125f052c90e326dcf6f380df94f6-Abstract-Datasets_and_Benchmarks_Track.html.

- [47] B. Sarkar, A. Shih, and D. Sadigh. Diverse conventions for human-AI collaboration. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=MljeRycu9s>.
- [48] J. Schulman, P. Moritz, S. Levine, M. I. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *CoRR*, abs/1506.02438, 2015. URL <https://api.semanticscholar.org/CorpusID:3075448>.
- [49] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *ArXiv*, abs/1707.06347, 2017. URL <https://api.semanticscholar.org/CorpusID:28695052>.
- [50] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [51] P. Stone, G. Kaminka, S. Kraus, and J. Rosenschein. Ad Hoc Autonomous Agent Teams: Collaboration without Pre-Coordination. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 24, pages 1504–1509, July 2010. doi: 10.1609/aaai.v24i1.7529. URL <https://ojs.aaai.org/index.php/AAAI/article/view/7529>.
- [52] D. Strouse, K. McKee, M. Botvinick, E. Hughes, and R. Everett. Collaborating with Humans without Human Data. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. S. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 14502–14515, 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/797134c3e42371bb4979a462eb2f042a-Paper.pdf.
- [53] T. Taylor. Evolutionary innovations and where to find them: Routes to open-ended evolution in natural and artificial systems. *Artificial life*, 25(2):207–224, 2019.
- [54] A. Vezhnevets, Y. Wu, M. Eckstein, R. Leblond, and J. Z. Leibo. Options as responses: Grounding behavioural hierarchies in multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 9733–9742. PMLR, 2020.
- [55] V. Villin, T. K. Buening, and C. Dimitrakakis. A minimax approach to ad hoc teamwork, 2025. URL <https://arxiv.org/abs/2502.02377>.
- [56] C. Wang, A. Rahman, I. Durugkar, E. Liebman, and P. Stone. N-agent ad hoc teamwork. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=q7TxGUWlhD>.
- [57] R. Wang, J. Lehman, A. Rawal, J. Zhi, Y. Li, J. Clune, and K. Stanley. Enhanced poet: Open-ended reinforcement learning through unbounded invention of learning challenges and their solutions. In *International conference on machine learning*, pages 9940–9951. PMLR, 2020.
- [58] X. Wang, S. Zhang, W. Zhang, W. Dong, J. Chen, Y. Wen, and W. Zhang. ZSC-Eval: An Evaluation Toolkit and Benchmark for Multi-agent Zero-shot Coordination. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, Nov. 2024. URL <https://openreview.net/forum?id=9aXjIBLwKc#discussion>.
- [59] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256, 1992.
- [60] D. Xing, Q. Liu, Q. Zheng, and G. Pan. Learning with generated teammates to achieve type-free ad-hoc teamwork. In Z.-H. Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 472–478. International Joint Conferences on Artificial Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/66. URL <https://doi.org/10.24963/ijcai.2021/66>. Main Track.
- [61] K. Xue, Y. Wang, C. Guan, L. Yuan, H. Fu, Q. Fu, C. Qian, and Y. Yu. Heterogeneous Multiagent Zero-Shot Coordination by Coevolution. *IEEE Transactions on Evolutionary Computation*, 29(5):2229–2243, Oct. 2025. ISSN 1941-0026. doi: 10.1109/TEVC.2024.3485177. URL <https://ieeexplore.ieee.org/document/10730789>.

- 642 [62] X. Yan, J. Guo, X. Lou, J. Wang, H. Zhang, and Y. Du. An efficient end-to-end training approach
643 for zero-shot human-AI coordination. In *Thirty-seventh Conference on Neural Information*
644 *Processing Systems*, 2023. URL <https://openreview.net/forum?id=6ePswuXUwf>.
- 645 [63] C. Yu, A. Velu, E. Vinitzky, J. Gao, Y. Wang, A. Bayen, and Y. Wu. The
646 Surprising Effectiveness of PPO in Cooperative Multi-Agent Games. In *Advances*
647 *in Neural Information Processing Systems*, volume 35, pages 24611–24624, Dec.
648 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/hash/](https://proceedings.neurips.cc/paper_files/paper/2022/hash/9c1535a02f0ce079433344e14d910597-Abstract-Datasets_and_Benchmarks.html)
649 [9c1535a02f0ce079433344e14d910597-Abstract-Datasets_and_Benchmarks.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/9c1535a02f0ce079433344e14d910597-Abstract-Datasets_and_Benchmarks.html).
- 650 [64] L. Yuan, L. Li, Z. Zhang, F. Chen, T. Zhang, C. Guan, Y. Yu, and Z.-H. Zhou. Learning to
651 coordinate with anyone. In *Proceedings of the Fifth International Conference on Distributed*
652 *Artificial Intelligence*, pages 1–9, 2023.
- 653 [65] R. Zhao, J. Song, Y. Yuan, H. Hu, Y. Gao, Y. Wu, Z. Sun, and W. Yang. Maximum entropy
654 population-based training for zero-shot human-ai coordination. In *Proceedings of the AAAI*
655 *Conference on Artificial Intelligence*, volume 37, pages 6145–6153, 2023.
- 656 [66] L. Zintgraf, S. Devlin, K. Ciosek, S. Whiteson, and K. Hofmann. Deep interactive bayesian
657 reinforcement learning via meta-learning. *arXiv preprint arXiv:2101.03864*, 2021.

658 Appendix

659 Please find an anonymized version of the code for this paper at <https://anonymous.4open.science/r/rotate/>.

661 A Algorithms ROTATE Pseudocode

662 ~~Open-Ended Ad Hoc Teamwork Environment, Env. Total of training iterations, T^{iter} . Initial ego~~
 663 ~~agent policy parameters, θ^{ego} . $B_{\pi} \leftarrow \langle \rangle$ $B_{\pi}^{\text{new}} \leftarrow \text{TeammateGenerator}(\text{Env}, \theta^{\text{ego}}, B_{\pi})$ $\theta^{\text{ego}} \leftarrow$~~
 664 ~~$\text{EgoUpdate}(\text{Env}, \theta^{\text{ego}}, B_{\pi}^{\text{new}})$ $B_{\pi} \leftarrow B_{\pi}^{\text{new}}$ **Return** θ^{ego}~~

665 A.1 Framework for Open Ended Ad Hoc Teamwork

666 Section 5 described an open-ended training framework for training an ego agent that can effectively
 667 collaborate with previously unseen teammates. ~~We further detail this general open-ended framework~~
 668 ~~, which is detailed~~ in Algorithm 1. ~~Here~~ In Line 3, a **TeammateGenerator** function determines a
 669 buffer of teammate policy parameters, B_{π}^{new} . The teammate generator function considers the ego
 670 agent’s current policy parameters, θ^{ego} , and the previous buffer of teammate policy parameters, B_{π}^{new} .
 671 Ideally, the teammate generation function generates and samples teammates that induce learning
 672 challenges to π^{ego} . In Line 4, an **EgoUpdate** function specifies a procedure that updates the ego
 673 agent’s policy parameters based on the B_{π}^{new} designed by the teammate generator. Pseudocode for
 674 ROTATE, which follows the open-ended framework specified by Algorithm 1, is presented in the
 675 following section.

676 A.2 ROTATE Algorithm

677 ROTATE’s teammate generation algorithm is detailed in Algorithm 2. As described in Section 6.1,
 678 this teammate generation algorithm jointly trains the parameters of a teammate policy and an estimate
 679 of its best response (BR) policy, based on a provided ego agent policy. The parameters of the teammate
 680 and BR policies, θ^{-i} and θ^{BR} , are initialized in Line 1. The parameters of the BR critic network, σ^{BR} ,
 681 are initialized in Line 2, while those for the teammate, $\sigma^{-i, \text{BR}}$ and $\sigma^{-i, \text{ego}}$, are initialized in Line 3.
 682 Note that the teammate maintains two critics, for separately estimating returns when interacting with
 683 the BR and ego agent policies.

684 The training of the teammate and BR policies is based on the SP, XP, XSP, and SXP interaction data
 685 gathered in Lines 5 to 8, which we previously motivated and described in Section 6.1. Recall that
 686 an SXP interaction require resetting an environment to start from an available XP state, and an XSP
 687 interaction analogously requires resetting to an SP state. Since resetting from all available XP states
 688 for SXP interaction is impractical, ROTATE *samples* from XP states to obtain start states for SXP
 689 interactions (and similarly for XSP). Experiences from SP, XP, SXP, and XSP interaction are stored in
 690 buffers $D_{\text{SP}}, D_{\text{XP}}, D_{\text{SXP}}, D_{\text{XSP}}$ in the form of a collection of tuples, $D = \langle (s_k, a_k, r_k, s'_k) \rangle_{k=1}^{|D|}$. Lines
 691 12 to 22 of Algorithm 2 then highlight how we use the stored experiences to compute loss functions
 692 that the trained models optimize.

693 Lines 12 and 13 describe how the teammate and BR policies are trained to mutually maximize
 694 returns when interacting with each other during SP and SXP interactions. Both lines call the
 695 **POL_LOSS_ADV_TARG** function, which receives $(\theta, \theta_{\text{old}}, \sigma_{\text{old}}, D, \epsilon)$ as input to evaluate the
 696 following, standard PPO-clip loss function that encourages return maximization and sufficient
 697 exploration:

$$\mathbb{E}_{(s, a, r, s') \in D} \left[\underbrace{-\min \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} A, \text{clip} \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A \right)}_{\text{PPO Clip Loss}} + \underbrace{\pi_{\theta}(a|s) \log(\pi_{\theta}(a|s))}_{\text{Entropy Loss}} \right],$$

698 where A denotes the *advantage* function. Our implementation of ROTATE uses an estimate of the
 699 advantage function obtained via the Generalized Advantage Estimation (GAE) algorithm [48], $A_{\sigma_{\text{old}}}^{\text{GAE}}$.
 700 Meanwhile, Line 14 shows how the teammate policy is trained to maximize the ego agent’s *regret*
 701 based on experiences from XP and SP interactions. The **POL_LOSS_REG_TARG** function that

Algorithm 2 ROTATE TeammateGenerator Function

Require:

Environment, Env.
 Ego agent policy, $\pi_{\theta^{\text{ego}}}$.
 Current teammate policy parameter buffer, B_π .
 Number of updates, N_{updates} .
 PPO clipping parameter, ϵ .
 PPO update epochs, N_{epochs} .

- 1: $\theta^{-i}, \theta^{\text{BR}} \leftarrow \text{RandomInit}(\pi), \text{RandomInit}(\pi)$
- 2: $\sigma^{\text{BR}} \leftarrow \text{RandomInit}(V)$
- 3: $\sigma^{-i, \text{BR}}, \sigma^{-i, \text{ego}} \leftarrow \text{RandomInit}(V), \text{RandomInit}(V) \quad \triangleright \text{Init teammate and BR parameters}$
- 4: **for** $t_{\text{update}} = 1, 2, \dots, N_{\text{updates}}$ **do**
- 5: $D_{\text{SP}}, D_{\text{XP}} \leftarrow \text{Interact}(\pi_{\theta^{\text{BR}}}, \pi_{\theta^{-i}}, p_0^{\text{Env}}), \text{Interact}(\pi_{\theta^{\text{ego}}}, \pi_{\theta^{-i}}, p_0^{\text{Env}})$
- 6: $s_{\text{XP}}, s_{\text{SP}} \leftarrow \text{SampleStates}(D_{\text{XP}}), \text{SampleStates}(D_{\text{SP}}) \quad \triangleright \text{Sample XP states}$
- 7: $D_{\text{SXP}} \leftarrow \text{Interact}(\pi_{\theta^{\text{BR}}}, \pi_{\theta^{-i}}, \mathcal{U}(s_{\text{XP}}))$
- 8: $D_{\text{XSP}} \leftarrow \text{Interact}(\pi_{\theta^{\text{ego}}}, \pi_{\theta^{-i}}, \mathcal{U}(s_{\text{SP}})) \quad \triangleright \text{Gather SP, XP, SXP, and XSP data}$
- 9: $\theta_{\text{old}}^{\text{BR}}, \theta_{\text{old}}^{-i} \leftarrow \theta^{\text{BR}}, \theta^{-i}$
- 10: $\sigma_{\text{old}}^{\text{BR}}, \sigma_{\text{old}}^{-i, \text{BR}}, \sigma_{\text{old}}^{-i, \text{ego}} \leftarrow \sigma^{\text{BR}}, \sigma^{-i, \text{BR}}, \sigma^{-i, \text{ego}} \quad \triangleright \text{Store old model parameters.}$
- 11: **for** $k_{\text{update}} = 1, 2, \dots, N_{\text{epochs}}$ **do**
- 12: $L_{\text{ppo-clip}}(\theta^{\text{BR}}) \leftarrow \text{POL_LOSS_ADV_TARG}(\theta_{\text{old}}^{\text{BR}}, \sigma_{\text{old}}^{\text{BR}}, D_{\text{SP}} \cup D_{\text{SXP}}, \epsilon)$
- 13: $L_{\text{ppo-clip}}(\theta^{-i}) \leftarrow \text{POL_LOSS_ADV_TARG}(\theta_{\text{old}}^{-i}, \sigma_{\text{old}}^{-i, \text{BR}}, D_{\text{SXP}}, \epsilon)$
- 14: $L_{\text{reg}}(\theta^{-i}) \leftarrow \text{POL_LOSS_REG_TARG}(\theta^{-i}, \theta_{\text{old}}^{-i}, \sigma_{\text{old}}^{-i, \text{BR}}, \sigma_{\text{old}}^{-i, \text{ego}}, D_{\text{SP}} \cup D_{\text{XP}}, \epsilon)$
- 15: $L_V(\sigma^{\text{BR}}) \leftarrow \text{VAL_LOSS}(\sigma_{\text{old}}^{\text{BR}}, D_{\text{SP}} \cup D_{\text{SXP}})$
- 16: $L_V(\sigma^{-i, \text{BR}}) \leftarrow \text{VAL_LOSS}(\sigma_{\text{old}}^{-i, \text{BR}}, D_{\text{SP}} \cup D_{\text{SXP}})$
- 17: $L_V(\sigma^{-i, \text{ego}}) \leftarrow \text{VAL_LOSS}(\sigma_{\text{old}}^{-i, \text{ego}}, D_{\text{XP}} \cup D_{\text{XSP}})$
- 18: $\theta^{\text{BR}} \leftarrow \text{GradDesc}(\theta^{\text{BR}}, \nabla_{\theta^{\text{BR}}} L_{\text{ppo-clip}}(\theta^{\text{BR}}))$
- 19: $\theta^{-i} \leftarrow \text{GradDesc}(\theta^{-i}, \nabla_{\theta^{-i}} (L_{\text{ppo-clip}}(\theta^{-i}) + L_{\text{reg}}(\theta^{-i}))) \quad \triangleright \text{Update policies}$
- 20: $\sigma^{\text{BR}} \leftarrow \text{GradDesc}(\sigma^{\text{BR}}, \nabla_{\sigma^{\text{BR}}} L_V(\sigma^{\text{BR}}))$
- 21: $\sigma^{-i, \text{BR}} \leftarrow \text{GradDesc}(\sigma^{-i, \text{BR}}, \nabla_{\sigma^{-i, \text{BR}}} L_V(\sigma^{-i, \text{BR}}))$
- 22: $\sigma^{-i, \text{ego}} \leftarrow \text{GradDesc}(\sigma^{-i, \text{ego}}, \nabla_{\sigma^{-i, \text{ego}}} L_V(\sigma^{-i, \text{ego}})) \quad \triangleright \text{Update critics.}$
- 23: **end for**
- 24: **end for**
- 25: $B_\pi \leftarrow B_\pi \cup \langle \theta^{-i} \rangle \quad \triangleright \text{Add generated teammate policy parameter}$
- 26: **Return** B_π

Algorithm 3 ROTATE EgoUpdate Function

Require:

Environment, Env.
 Ego agent policy parameters, θ^{ego} .
 Current teammate policy parameter buffer, B_π .
 Number of updates, N_{updates} .
 PPO clipping parameter, ϵ .
 PPO update epochs, N_{epochs}

```

1:  $\sigma^{\text{ego}} \leftarrow \text{Init}(V)$  ▷ Init params of the critic networks of  $\pi^{\text{ego}}$ 
2: for  $t_{\text{update}} = 1, 2, \dots, N_{\text{updates}}$  do
3:    $\theta^{-i} \sim \mathcal{U}(B_\pi)$  ▷ Sample teammate parameters uniformly
4:    $D \leftarrow \text{Interact}(\pi_{\theta^{-i}}, \pi_{\theta^{\text{ego}}}, p_0^{\text{Env}})$ 
5:    $\theta_{\text{old}}^{\text{ego}}, \sigma_{\text{old}}^{\text{ego}} \leftarrow \theta^{\text{ego}}, \sigma^{\text{ego}}$ 
6:   for  $k_{\text{update}} \in \{1, 2, \dots, N_{\text{epochs}}\}$  do
7:      $L_\pi(\theta^{\text{ego}}) \leftarrow \text{EGO\_POL\_LOSS}(\theta^{\text{ego}}, \theta_{\text{old}}^{\text{ego}}, \sigma_{\text{old}}^{\text{ego}}, D, \epsilon)$  ▷ Compute policy loss
8:      $L_V(\sigma^{\text{ego}}) \leftarrow \text{EGO\_VAL\_LOSS}(\sigma^{\text{ego}}, \sigma_{\text{old}}^{\text{ego}}, D, \epsilon)$  ▷ Compute critic loss
9:      $\theta^{\text{ego}} \leftarrow \text{GradDesc}(\theta^{\text{ego}}, \nabla_{\theta^{\text{ego}}} L_\pi(\theta^{\text{ego}}))$  ▷ Update policy
10:     $\sigma^{\text{ego}} \leftarrow \text{GradDesc}(\sigma^{\text{ego}}, \nabla_{\sigma^{\text{ego}}} L_V(\sigma^{\text{ego}}))$  ▷ Update critic
11:   end for
12: end for
13: Return  $\theta^{\text{ego}}$ 

```

702 computes a loss function that encourages the maximization of regret is generally the same as the
 703 **POL_LOSS_ADV_TARG** function except for its replacement of the advantage function, A , with a
 704 regret-based target function. The regret-based target function is defined differently but symmetrically
 705 for SP and XP states. We describe the target function for regret from **XP** states below, and refer the
 706 reader to the code for the target function for regret from SP states.

$$A_{\text{reg}} = \underbrace{V_{\sigma_{\text{old}}^{-i, \text{BR}}}(s)}_{\approx V(s|\pi^{-i}, \text{BR}(\pi^{-i}))} - \underbrace{(r + \gamma V_{\sigma_{\text{old}}^{-i, \text{ego}}}(s'))}_{\approx V(s|\pi^{-i}, \pi^{\text{ego}})}. \quad (11)$$

707 Rather than optimizing a regret function that requires explicitly computing the return-to-go,
 708 **POL_LOSS_REG_TARG** estimates the XP return via a 1-step bootstrapped return using the
 709 teammate critic parameterized by $\sigma^{-i, \text{ego}}$. Similarly, the SP return is estimated using the teammate
 710 critic network parameterized by $\sigma^{-i, \text{BR}}$. This results in a regret optimization method that uses the
 711 log-derivative trick to optimize objective functions [59, 21]. The ROTATE regret estimation method
 712 and alternative approaches to maximize regret are further discussed in App. [C.4](#).

713 Lines 15 to 17 then detail how we train critic networks that measure returns from the interaction
 714 between the generated teammate policy and its best response or ego agent policy. We specifically call
 715 the **VAL_LOSS** function that receives $(\sigma, \sigma_{\text{old}}, D)$ to compute the standard mean squared Bellman
 716 error (MSBE) loss, defined as:

$$\mathbb{E}_{(s, a, r', s') \in D} \left[(V_\sigma(s) - V_{\sigma_{\text{old}}}^{\text{targ}}(s))^2 \right], \quad (12)$$

717 where $V_{\sigma_{\text{old}}}^{\text{targ}}(s) := A_{\sigma_{\text{old}}}^{\text{GAE}} + V_{\sigma_{\text{old}}}(s)$ is the target value estimate.

718 The previously defined loss functions can be minimized using any gradient descent-based optimization
 719 technique, as we indicate in Lines 18 to 22. In practice, our implementation uses the ADAM
 720 optimization technique [26]. At the end of this teammate generation process, Lines 25 and 26 indicate
 721 how the generated teammate policy parameter is added to a storage buffer, which is subsequently
 722 uniformly sampled to provide teammate policies for ego agent training.

723 The ego agent policy’s training process proceeds according to Algorithm 3. Line 3 illustrates how
 724 ROTATE creates different teammate policies by uniformly sampling model parameters from the
 725 B_π resulting from the teammate generation process. Using the experience collaborating with the
 726 sampled policies outlined in Line 4, the ego agent’s policy parameters are updated to maximize

its returns via PPO in Line 7. The only difference between the **EGO_POL_LOSS** function and **POL_LOSS_ADV_TARG** function in Algorithm 2 is the input used to compute the loss function. Unlike in the **EGO_POL_LOSS** function, we assume that the input dataset, D , stores the historical sequence of observed states and executed actions, h , rather than states. Likewise, we assume that the only difference between the **VAL_LOSS** and **EGO_VAL_LOSS** function is that the latter stores the observation-action history rather than states (Line 8). Like recent AHT learning algorithms [66, 41, 40], π^{ego} and V^{ego} are conditioned on the ego agent’s observation-action history to facilitate an adaptive π^{ego} through an improved characterization of teammates’ policies. The history-conditioned ego architecture and other practical implementation details are described in App. F. Finally, the ego agent update function returns the updated ego agent policy parameters, which are provided as part of the inputs for the next call to ROTATE’s teammate generation function.

B Baselines Overview

The main paper compares ROTATE to five baselines: PAIRED, Minimax Return, FCP, BRDiv, and CoMeDi. Each baseline is briefly described below, followed by a discussion of the computational complexity of teammate generation baselines compared to ROTATE, and a discussion of the relationship of Mixed Play (MP) with per-state and per-trajectory regret. A discussion of implementation details can be found in App. F.

PAIRED [17]: A UED algorithm where a regret-maximizing “adversary” agent proposes environment variations that an allied antagonist achieves high returns on, but a protagonist agent receives low returns on. The algorithm is directly applicable to AHT by defining a teammate generator for the role of the adversary, a best response agent to the generated teammate for the role of the antagonist, and an ego agent for the role of the protagonist.

Minimax Return [37, 55]: A common baseline in the UED literature, with origins in robust reinforcement learning, where the objective is minimax return. Prior works in AHT have proposed generating a curriculum of teammates according to this objective. Translated to our open-ended learning setting, the teammate generator creates teammates that minimize the ego agent’s return, while the ego agent maximizes return.

Fictitious Co-Play [52]: A two-stage AHT algorithm where a pool of teammates is generated by running IPPO [63] with varying seeds, and saving multiple checkpoints to the pool. The ego agent is an IPPO agent that is trained against the pool.

BRDiv [42]: A two-stage AHT algorithm where a population of “confederate” and best-response agent pairs is generated, and an ego agent is trained against the confederates. BRDiv maintains a cross-play matrix containing the returns for all confederate and best-response pairs. The diagonal returns (self-play) are maximized, while the off-diagonal returns (cross-play) are minimized. BRDiv and LIPO [10] share a similar objective, where the main differences are: (1) If xp_weight denotes the weight on the XP return, then BRDiv requires that the coefficient on the SP return is always $1 + 2 * xp_weight$, and (2) LIPO introduces a secondary diversity metric based on mutual information, and (3) LIPO assumes that agents within a team (i.e., a confederate-BR pair) share parameters.

CoMeDi [47]: CoMeDi is a two-stage AHT algorithm. In the first stage, a population of teammates is generated, and in the second stage, an ego agent is trained against the teammate population. The teammate generation stage trains teammate policies one at a time, where the n th teammate policy is trained to maximize its SP return, minimize its XP return with the previously generated teammate (i.e. from among teammates $1, \dots, n - 1$) that it best collaborates with, and maximizes its “mixed-play” (MP) return. The relationship between the regret objectives described in Section 6 and MP is further discussed in App. [C.3](#).

B.1 Computational Complexity of ROTATE versus Teammate Generation Baselines

The computational complexity of ROTATE is compared with that of the teammate generation baselines, in terms of the population size and the number of objective updates. In the following, n denotes the population size, while T indicates the number of updates needed to train an individual

776 population member. The precise meaning of n and T might vary with the algorithm, but is made
777 clear in each description.

778 **FCP:** Let T denote the number of RL updates needed to train each IPPO team and let n denote the
779 number of teams trained by FCP. Then, the computational complexity of FCP is $O(nT)$.

780 **BRDiV/LIPO:** Both BRDiV [42] and LIPO [10] require sampling trajectories from each pair of
781 agents in the population, for each update. Thus, if the total number of updates is T and the population
782 size is n , then the algorithm’s time complexity is $O(n^2T)$. Due to the quadratic complexity in n ,
783 BRDiV and LIPO are typically run with smaller population sizes, with $n < 10$ for all non-matrix game
784 tasks in both original papers.

785 **CoMeDi:** Recall that CoMeDi trains population members one at a time, such that each agent is
786 distinct from the previously discovered teammates in the population. This necessitates performing
787 evaluation rollouts of the currently trained agent against all previously generated teammates at each
788 RL update step. Let T be the number of RL updates required to train the i th agent to convergence,
789 and let n denote the population size. Then CoMeDi’s time complexity is $O(n^2T)$ —making it scale
790 quadratically in n , similar to BRDiV and LIPO.

791 **ROTATE:** In ROTATE, a new teammate is trained to convergence for each iteration of open-ended
792 learning. Thus, the number of open-ended learning iterations is equal to the population size n , where
793 within each iteration, there are $O(T)$ RL updates performed. Therefore, the complexity of ROTATE
794 is $O(nT)$, meaning that our method scales linearly in the population size n .

795 C Supplemental Results Discussion

796 ~~This section presents various supplemental results. First~~Here, we describe how ROTATE satisfies
797 the definition of an open-ended system given by Hughes et al. [23]. Next, we compare and contrast
798 CoMeDi’s mixed-play mechanism in the context of to ROTATE’s per-state regret. ~~Second~~Finally,
799 we discuss alternative estimators for ROTATE per-state regret. ~~Third, we present experiments~~
800 ~~comparing to a variant with CoMeDi-style mixed-play return maximization.~~

801 C.1 Formalizing ROTATE as an Open-Ended System

802 C.1.1 Definition of Open-ended Systems

803 To formalize open-endedness, we employ the definition of open-ended systems provided by [23].
804 [23] consider a system S that produces a sequence of artifacts X_t , indexed by time t . An observer
805 O attempts to predict a new artifact X_T based on observations of past artifacts seen up until time t .

806 From the perspective of an observer, a system is *open-ended* if and only if the
807 sequence of artifacts it produces is both novel and a variant using the alternative
808 regret estimation strategy. ~~Fourth, we examine whether the population generated~~
809 ~~by ROTATE is useful for training an independent ego agent~~learnable. ~~Fifth, we~~
810 ~~present and describe radar charts breaking down the performance of on all six~~
811 ~~benchmark tasks presented in~~

812 A sequence of artifacts is novel if it becomes increasingly unpredictable with respect to an observer’s
813 model at any fixed time t . Formally:

$$\forall t, \forall T > t, \exists T' > T : \mathbb{E}[\ell(t, T')] > \mathbb{E}[\ell(t, T)]. \quad (13)$$

814 This is quantified by a loss metric, $\ell(t, T)$, where the first argument refers to the timestep that the
815 observer has observed up until (i.e., artifacts X_1, \dots, X_t), while the second argument refers to the
816 artifact that the observer must now predict (i.e., X_T).

817 Conversely, the system is learnable whenever conditioning on a longer history makes artifacts more
818 predictable:

$$\forall T, \forall t < T, \forall T' > t : \mathbb{E}[\ell(t', T)] < \mathbb{E}[\ell(t, T)]. \quad (14)$$

C.2 Open-endedness in ROTATE

ROTATE displays both novelty and learnability according to the definition above, as we show below. Let the ego agent be the observer, the sequence of generated teammates be the artifacts, and $\ell(t, T)$ be the cooperative regret of the ego agent trained up to iteration t , evaluated on the teammate generated at iteration T .

Novelty. At open-ended iteration $t + 1$, ROTATE seeks to generate a teammate that *challenges* the ego agent from iteration t by directly maximizing its cooperative regret. By construction, for $T' = T + 1$:

$$\mathbb{E}[\ell(t, T')] > \mathbb{E}[\ell(t, T)]. \quad (15)$$

Learnability. In the context of ROTATE, this condition states that for some *future unseen teammate from time T* , the cooperative regret of ROTATE’s ego agent from iteration t' should be lower than the cooperative regret of ROTATE’s ego agent from an earlier iteration t .

Intuitively, since all generated teammates are added to a population buffer that is uniformly sampled, ROTATE’s ego agent at future iterations should always be more capable than the ego agent at previous iterations. Empirically, we observe this in the learning curves presented in Fig. 8. The normalized return against the test set of unseen teammates increases as the iterations increase. Importantly, the main paper, and the learning curves for all variants of that are tested in this paper. Finally, we present a human proxy evaluation of ROTATE and CoMeDi on the Overcooked tasks. return is min-max normalized so that 1.0 represents the return achieved by a teammate with its best response, while 0.0 represents a known lower bound on returns. If the normalized return increases from iteration t to t' , then the cooperative regret decreases in the same interval. We provide the formal argument below.

C.2.1 Relationship Between Normalized Return and Cooperative Regret

Let η represent the lower return bound. BR is used as a shorthand for $\text{BR}(\pi^{-i})$. Note that $V(\pi^{-i}, \text{BR})$ is the upper bound on returns for teammate π^{-i} . We denote the cooperative regret (CR) and min-max normalized returns of an ego agent at open-ended iteration t with a fixed teammate π^{-i} as:

$$\text{CR}(t) := V(\pi^{-i}, \text{BR}) - V(\pi^{-i}, \pi_t^{\text{ego}}) \quad (16)$$

$$\text{NRet}(t) := \frac{V(\pi^{-i}, \pi_t^{\text{ego}}) - \eta}{V(\pi^{-i}, \text{BR}) - \eta} \quad (17)$$

Proposition 1. Consider the learning interval $[t, t']$ where $t < t'$ and $V(\pi^{-i}, \pi_t^{\text{ego}}) < V(\pi^{-i}, \pi_{t'}^{\text{ego}})$. Then, $\text{CR}(t) > \text{CR}(t')$ and $\text{NRet}(t) < \text{NRet}(t')$.

First, we transform the regret to guarantee each term is nonnegative by adding and subtracting η :

$$\text{CR}(t) = (V(\pi^{-i}, \text{BR}) - \eta) - (V(\pi^{-i}, \pi_t^{\text{ego}}) - \eta) \quad (18)$$

Now, let c denote the positively transformed return of the teammate with its best response, which is constant as π^{ego} varies:

$$c := V(\pi^{-i}, \text{BR}) - \eta. \quad (19)$$

Further, let $r(t)$ denote the positively transformed returns of the ego agent at iteration t :

$$r(t) = V(\pi^{-i}, \pi_t^{\text{ego}}) - \eta. \quad (20)$$

Thus, the regret and normalized return of an ego agent at iteration t can be rewritten as:

$$\text{CR}(t) = c - r(t) \quad (21)$$

$$\text{NRet}(t) = \frac{r(t)}{c} \quad (22)$$

By assumption, over $[t, t']$, $r(t)$ increases monotonically. We analyze the above as linear equations varying in $r(t)$. Since $r(t) > 0$ and $c > 0$, the normalized return $\text{NRet}(t)$ increases in the interval $[t, t']$, while the regret $\text{Reg}(t)$ decreases in the same interval.

C.3 Discussion of Relationship To CoMeDi and Mixed Play

As previously described in App. B, CoMeDi [47] is a two-stage teammate generation AHT algorithm, whose teammate generation process trains one teammate per iteration, with an objective that encourages the new teammates to be distinct from previously discovered teammates.

CoMeDi adds trained teammates policies to a teammate policy buffer, Π^{train} . Each iteration begins by identifying the teammate policy that is most *compatible* with the currently trained teammate π^{-i} , out of all previously generated policies:

$$\pi^{\text{comp}} = \underset{\pi^{-j} \in \Pi^{\text{train}}}{\text{argmax}} \mathbb{E}_{s \sim p_0} [V(s | \pi^{-i}, \pi^{-j})]. \quad (23)$$

The new teammate policy π^{-i} is trained with an objective that improves the per-trajectory regret objective (Eq. 8) by adding a term that maximizes the returns from states gathered in *mixed-play*, which we describe below.

Let mixed-play starting states be sampled from states visited when π^{-i} interacts with the *mixed* policy, that uniformly samples actions from π^{comp} and $\text{BR}(\pi^{-i})$ at each timestep:

$$p_{\text{MSTART}} := d \left(\pi^{-i}, \frac{1}{2} \pi^{\text{comp}} + \frac{1}{2} \text{BR}(\pi^{-i}); p_0 \right). \quad (24)$$

From these starting states, CoMeDi then gathers mixed-play interaction data, where π^{-i} interacts with $\text{BR}(\pi^{-i})$. The resulting mixed-play state visitation is then expressed as:

$$p_{\text{MP}} := d \left(\pi^{-i}, \text{BR}(\pi^{-i}); p_{\text{MSTART}} \right). \quad (25)$$

The complete objective that Sarkar et al. [47] optimizes to train a collection of diverse teammates is then defined as:

$$\max_{\pi} \left(\mathbb{E}_{s_0 \sim p_0} [\text{CR}(\pi^{\text{comp}}, \pi^{-i}, s_0)] + \underbrace{\mathbb{E}_{s \sim p_{\text{MP}}} [V(s | \pi, \text{BR}(\pi))]}_{\text{mixed-play return maximization}} \right). \quad (26)$$

CoMeDi [47] optimizes this objective to discourage π^{-i} from learning poor actions for collaborations outside of p_{SP} . This is because π^{-i} is now also trained to maximize returns in states visited during mixed-play, which resembles some states encountered while cooperating with π^{comp} . Discerning whether a state is likely encountered while interacting with π^{comp} and consequently choosing to sabotage collaboration will no longer be an optimal policy to maximize Expr. 26.

Despite the importance of using p_{MSTART} as a starting state for data collection being questionable, we take inspiration from CoMeDi’s maximization of $V(s | \pi, \text{BR}(\pi))$ outside of states from p_{SP} . We argue that maximizing $V(s | \pi^{-i}, \text{BR}(\pi^{-i}))$ is a key component towards making π^{-i} act in good faith by always choosing actions yielding optimal collective returns assuming $\text{BR}(\pi^{-i})$ is substituted as the partner policy. Unlike CoMeDi, ROTATE maximizes $V(s | \pi^{-i}, \text{BR}(\pi^{-i}))$ on trajectories gathered from a starting state from p_{XP} (i.e., SXP states) instead of p_{MSTART} , which results in the second term of Expr. 10. We formulate this objective to encourage π^{-i} to act in good faith in states sampled from p_{XP} , which is visited while π^{-i} interacts with π^{ego} . Since π^{-i} is not sabotaging π^{ego} by selecting actions that make collaboration impossible in p_{XP} , the ego policy learning process becomes less challenging. We conjecture that this leads to π^{ego} with better performances as indicated in Figure 3.

While Figure 3 compares ROTATE with CoMeDi, Figure 6a compares ROTATE with a modified CoMeDi approach that now follows the open-ended training framework described in Algorithm 1. In this modified version of CoMeDi, we train a newly generated teammate policy to maximize Eq. 26 while substituting π^{comp} with the trained π^{ego} . Rather than promoting meaningful differences with previously generated teammate policies, this creates a teammate policy that maximizes the ego agent policy’s per-trajectory regret while mitigating self-sabotage. This version of CoMeDi’s teammate generation objective within the ROTATE open-ended framework is visualized in Figure 5.

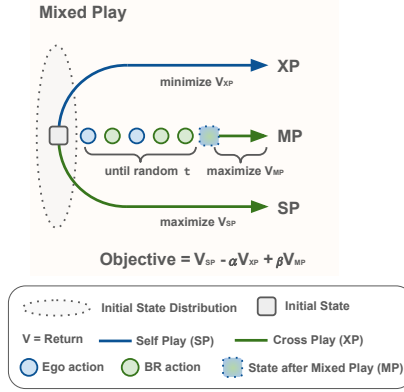


Figure 5: CoMeDi-style mixed-play objective for teammate generation, in the context of open-ended AHT.

893 C.4 Alternatives Estimators for Per-State Regret

894 This section discusses the approach employed by ROTATE in Algorithm 2 to estimate the per-state
 895 regret objective under a specific distribution, as well as an alternative estimation method. Experiments
 896 comparing the two approaches are also presented and discussed.

897 Recall that the per-state regret under states sampled from a distribution D is defined as:

$$\mathbb{E}_{s \sim D} [\text{CR}(\pi^{\text{ego}}, \pi^{-i}, s)] = \mathbb{E}_{s \sim D} [V(s|\pi^{-i}, \text{BR}(\pi^{-i})) - V(s|\pi^{-i}, \pi^{\text{ego}})] \quad (27)$$

$$= \underbrace{\mathbb{E}_{s \sim D} [V(s|\pi^{-i}, \text{BR}(\pi^{-i}))]}_{\text{SP return}} - \underbrace{\mathbb{E}_{s \sim D} [V(s|\pi^{-i}, \pi^{\text{ego}})]}_{\text{XP return}}. \quad (28)$$

898 In practice, we can use the policy gradient method to maximize regret by estimating the self-play
 899 returns and cross-play returns in Eq. 28 using the n -step return, Monte Carlo-based return-to-go
 900 estimate, or generally any variant of the advantage function estimator. The choice of return estimates
 901 affects the result of our teammate generation process through the bias-variance tradeoff when
 902 estimating regret. Combined with the potentially different choices of D , we can design different
 903 variants of ROTATE based on how regret is estimated.

904 **ROTATE Per-State Regret:** Line 14 in Algorithm 2 and Eq. 11 outline how ROTATE maximizes
 905 per-state regret in states visited during XP interaction (denoted by p_{XP}), where SP and XP returns
 906 are estimated via a trained critic and a 1-step return estimate, respectively. As a reminder, ROTATE
 907 employs the following regret target function to train the regret-maximizing teammate policy on XP
 908 states, with an analogously defined target function for SP states:

$$\mathbb{E}_{s \sim p_{\text{XP}}} \left[\underbrace{V_{\sigma^{-i}, \text{BR}}(s)}_{\text{SP return estimate}} - \underbrace{(r + \gamma V_{\sigma^{-i}, \text{ego}}(s'))}_{\text{XP return estimate}} \right]. \quad (29)$$

909 We maximize regret in states sampled from p_{XP} and p_{SP} to encourage the design of teammate policies
 910 that provide a learning challenge while also acting in good faith, thereby maximizing cooperative
 911 returns assuming interactions with its best-response policy, while interacting with the ego agent's
 912 policy. Our discussion here focuses on computing Eq. 29 for brevity. However, a similar approach
 913 can be used to train a critic network to estimate regret in SP states accurately. The only difference
 914 lies in the use of states sampled from p_{SP} and p_{XSP} for training the critic network.

915 Despite potentially providing biased estimates, training a value function to estimate self-play
 916 returns can reduce the variance caused by environment stochasticity, compared to a Monte Carlo return-to-go
 917 estimate.

918 The critic network estimating teammate-BR returns, $V_{\sigma^{-i}, \text{BR}}(s)$, is trained on interactions initialized
 919 from XP, as shown in Line 16 of Algorithm 2. This enables the teammate-BR critic network to
 920 accurately estimate SP returns from p_{XP} states. Meanwhile, a 1-step estimate of XP returns is made
 921 possible by storage of rewards experienced during XP interactions (Line 5 of Algorithm 2) and the

training of a value function to estimate XP returns (Line 17 of Algorithm 2). Utilizing a 1-step estimate produces lower variance than using a Monte Carlo-based return-to-go estimate, while also yielding less bias than predicting returns solely based on the trained critic network’s value.

Estimating Per-State Regret via Monte Carlo Returns: An alternative approach for estimating is to use a Monte Carlo-based return-to-go estimate for both SP and XP return estimates. Assuming that both interaction starts from states encountered during XP interaction, the policy updates under this alternative approach maximize the following target function:

$$\mathbb{E}_{s_t \sim 0.5p_{\text{XP}} + 0.5p_{\text{SP}}} \left[\underbrace{\mathbb{E}_{a_{t'} \sim [\text{BR}(\pi^{-i}), \pi^{-i}], P} \left[\sum_{t'=t}^{\infty} \gamma^{t'} r_{t'} \middle| s_t \right]}_{\text{SP return estimate}} - \underbrace{\mathbb{E}_{a_{t'} \sim [\pi^{\text{ego}}, \pi^{-i}], P} \left[\sum_{l=0}^{\infty} \gamma^{t'} r_{t'} \middle| s_t \right]}_{\text{XP return estimate}} \right]. \quad (30)$$

We refer to this as the *Monte Carlo* per-state regret. However, starting both SP and XP interactions from all states visited in XP can be computationally prohibitive. More importantly, the Monte Carlo-based return-to-go estimates of SP and XP returns have high variance, especially when the environment transition function and the trained policies are highly stochastic.

Estimating Per-State Regret via Generalized Advantage Estimators: A final approach for estimating Eq. 27 is to substitute both return-to-go estimates in Expr. 30 with a generalized advantage estimator [48] based on SP and XP interactions. This results in the maximization of the following target function during the teammate policy updates:

$$\mathbb{E}_{s_t \sim 0.5p_{\text{XP}} + 0.5p_{\text{SP}}} \left[\underbrace{\mathbb{E}_{a_{t'} \sim [\text{BR}(\pi^{-i}), \pi^{-i}], P} \left[\underbrace{\sum_{t'=t}^{\infty} (\gamma\lambda)^{t'} \delta_{t'}^{-i, \text{BR}}}_{\text{GAE}} \middle| s_0 \right]}_{\text{SP return estimate}} - \underbrace{\mathbb{E}_{a_{t'} \sim [\pi^{\text{ego}}, \pi^{-i}], P} \left[\underbrace{\sum_{t'=t}^{\infty} (\gamma\lambda)^{t'} \delta_{t'}^{-i, \text{ego}}}_{\text{GAE}} \middle| s_0 \right]}_{\text{XP return estimate}} \right], \quad (31)$$

where we define $\delta_t^{-i, \text{BR}}$ and $\delta_t^{-i, \text{ego}}$ as:

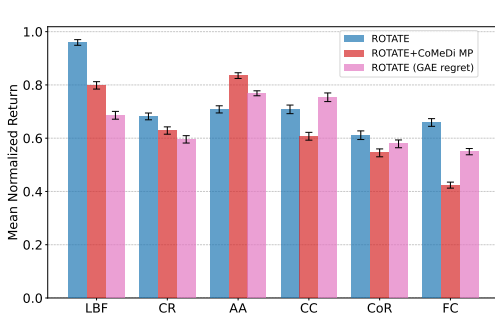
$$\begin{aligned} \delta_t^{-i, \text{BR}} &= r_t + \gamma V_{\sigma^{-i, \text{BR}}}(s_{t+1}) - V_{\sigma^{-i, \text{BR}}}(s_t), \\ \delta_t^{-i, \text{ego}} &= r_t + \gamma V_{\sigma^{-i, \text{ego}}}(s_{t+1}) - V_{\sigma^{-i, \text{ego}}}(s_t). \end{aligned}$$

We refer to an instance of the ROTATE algorithm that maximizes regret using this target function as ROTATE with *GAE per-state regret*. In practice, we collect data for SP GAE maximization and XP GAE minimization by first independently sampling two collections of states from D_{SXP} and D_{XP} respectively. Next, the states sampled from D_{SXP} are used to maximize the GAE from SXP interactions, while states sampled from D_{XP} are utilized to minimize the GAE from XP interactions. The γ and λ parameters used during the computation of the generalized advantage estimator are mechanisms to regulate the bias and variance of the regret estimation [48], effectively providing a different bias-variance tradeoff compared to the previously mentioned methods.

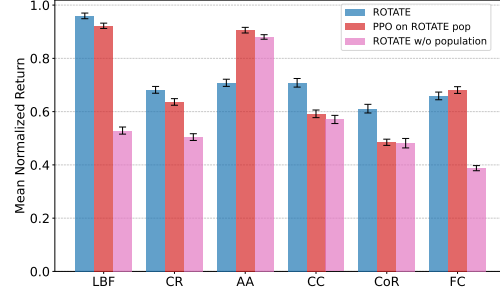
D Supplemental Experiments and Results

This section provides additional experiments examining the ROTATE teammate generation objective and population buffer in further detail. It also provides supplemental material supporting the core results. Specifically:

- We compare ROTATE to variants of ROTATE that uses CoMeDi-style mixed-play return maximization (App. C.3) and a GAE-based regret estimation strategy (App. C.4).
- We examine whether the population generated by ROTATE is useful for training an independent ego agent.



(a) ROTATE vs ROTATE with CoMeDi’s mixed-play (MP) objective and ROTATE with GAE regret.



(b) ROTATE compared to an independently trained ego agent on ROTATE’s population and an ablation of ROTATE without the population.

Figure 6: The mean normalized returns of ROTATE and various ablations designed to evaluate the effectiveness of ROTATE’s regret-based teammate generation objective and population-based ego agent training procedure.

- [We present radar charts breaking down the performance of ROTATE on all six benchmark tasks presented in the main paper, and the learning curves for all variants of ROTATE that are tested in this paper](#)
- [We provide a human proxy evaluation of ROTATE and CoMeDi on the Overcooked tasks.](#)

D.1 Experimental Comparisons of ROTATE with Alternative Teammate Generation Objectives

Figure 6a compares the version of ROTATE presented in the main paper and Algorithm 2, to ROTATE with GAE per-state regret, and a version of ROTATE where expected returns are maximized in states sampled from p_{MP} rather than p_{SXP} , which resembles the mixed-play objective of CoMeDi [47]. We do not implement the Monte Carlo per-state regret estimation approach described above, as it is impractical and unlikely to yield better results than using value functions to estimate regret. ROTATE and ROTATE with GAE regret yield mixed results as neither approach consistently beats the other in all environments. We suspect this is caused by the policy gradient’s different bias and variance levels when estimating regret using these two methods. Meanwhile, ROTATE’s maximization of returns in states from p_{SXP} leads to higher normalized returns than maximizing CoMeDi’s mixed-play objective in all environments except for Overcooked’s Asymmetric Advantages (AA) setting. Following the difference in starting states of trajectories for which these two maximize self-play returns, we conjecture that this is because ROTATE empirically teammate policies with good faith in states from p_{XP} while the CoMeDi-like approach imposes the same thing in states from p_{MSTART} . Imposing good faith within policies in p_{XP} is likely more important for training an ego agent that initially interacts with π^{-i} during training by visiting states from p_{XP} .

D.2 Training an Independent Ego Agent on the ROTATE Population

Two-stage AHT algorithms first generate a population of teammates, and next train an ego agent against the population. Although ROTATE’s teammate generation mechanism relies on the learning process of a particular ego agent, here, we investigate whether the population generated by ROTATE is useful for training independently generated ego agents. Fig. 6b compares the mean evaluation returns of the ROTATE ego agent against the mean evaluation returns of an independently trained ego agent that was trained using the same configuration as ROTATE. In 3/6 tasks, the ROTATE ego agent outperforms the trained ego agent, while in two tasks, the two ego agents perform similarly (LBF and FC). This result suggests that the ROTATE population is a useful population of teammates even independent of the particular ego agent generated. The strong performance of the independently trained ego agent is unsurprising given that it has two advantages over the ROTATE ego agent. First, the independently trained ego agent faces a stationary distribution of training teammates compared to ROTATE, which faces a nonstationary distribution caused by the population growing over learning iterations. Second, the independently trained ego agent interacts with all teammates uniformly throughout training, while the ROTATE ego agent only trains against earlier teammates for more iterations than later teammates.

D.3 ROTATE vs Baselines—Radar Charts

We break down the performance of ROTATE and all baseline methods by individual evaluation teammate policies as radar charts in Fig. 7. The radar charts show that ROTATE achieves higher performance across a larger number and variety of evaluation teammates than baselines. The best baseline, CoMeDi, achieves unusually high returns with the heuristic-based evaluation teammates on LBF, CR, and CC. We hypothesize that this trend occurs because CoMeDi explicitly optimizes for novel conventions that do not match existing conventions. However, on these tasks, CoMeDi does not perform as well as BRDiv teammates, which are trained to maximize the adversarial diversity objective. The radar charts also show that the second-best baseline, FCP, is strong specifically against IPPO teammates and relatively weaker on heuristics and BRDiv teammates, especially in CR and CC. As mentioned in the main paper, we attribute FCP’s relative strength on IPPO evaluation teammates to the fact that the IPPO evaluation teammates are closer to the training teammate distribution constructed by FCP. While FCP is not especially strong against the “IPPO pass” agents in CC, these agents were trained via reward shaping to solve the task by passing onions across the counter rather than navigating around the counter, which is the policy found by IPPO without reward shaping (denoted as “IPPO CC” in the figures).

D.4 Learning Curves

Figure 8 shows learning curves for ROTATE and all ROTATE variations tested in this paper, where the x -axis is the open-ended learning iteration, and the y -axis corresponds to the mean evaluation return. On 4/6 tasks (LBF, CR, CC, and FC), ROTATE has better sample efficiency than variants. On 3/6 tasks (LBF, CR, and FC), ROTATE dominates variants at almost all points in learning.

D.5 Human Proxy Evaluations

The results in the main paper show that ROTATE is better able to generalize to unseen partners compared to baseline methods. Here, we evaluate the ability of ROTATE to generalize to *human* partners, compared to the best baseline, CoMeDi. The evaluation is conducted with human proxy teammates generated by behavior cloning on the human gameplay dataset published by [9]. Table 2 shows that ROTATE achieves higher returns in coordination with the human proxy compared to CoMeDi, across all five Overcooked layouts.

| | CR | AA | CC | CoR |
|--------|--|--|--|--|
| ROTATE | 0.810.85 (0.78, 0.84, 0.82, 0.87) | 0.660.68 (0.63, 0.69, 0.65, 0.70) | 1.090.92 (1.05, 1.14, 0.88, 0.95) | 0.730.66 (0.70, 0.63, 0.58, 0.69) |
| CoMeDi | 0.75 (0.71, 0.78) | 0.57 (0.54, 0.59) | 0.89 (0.83, 0.95) | 0.56 (0.54, 0.58) |

Table 2: ROTATE outperforms CoMeDi with the human proxy teammate on all Overcooked layouts. Normalized returns and bootstrapped 95% CI’s are shown, where the normalization is performed using the human proxy agent’s self-play returns.

E Experimental Tasks

Experiments in the main paper are conducted on Jax re-implementations of Level-Based Foraging (LBF) [2, 6], and five tasks from the Overcooked suite—~~Cramped suite~~ (Cramped Room (CR), Asymmetric Advantages (AA), Counter Circuit (CC), Coordination Ring (CoR), and Forced Coordination (FC) [9, 46]. ~~Each task is described in [9, 46] and a sabotage matrix game. The sabotage matrix game is already fully described in the main paper. The other tasks are described below.~~

Level-Based Foraging (LBF) Originally introduced by Albrecht and Ramamoorthy [2], Level-Based Foraging is a mixed cooperative-competitive logistics problem where N players interact within a rectangular grid world to obtain k foods. All players and foods have a positive integer *level*, where groups of one to four players may only *load* (collect) a food if the sum of player levels is greater than the food’s level. A food’s level is configured so that it is always possible to load it.

We use the Jax re-implementation of LBF by Bonnet et al. [6], which was based on the implementation by Christianos et al. [13]. The implementation permits the user to specify the number of players,

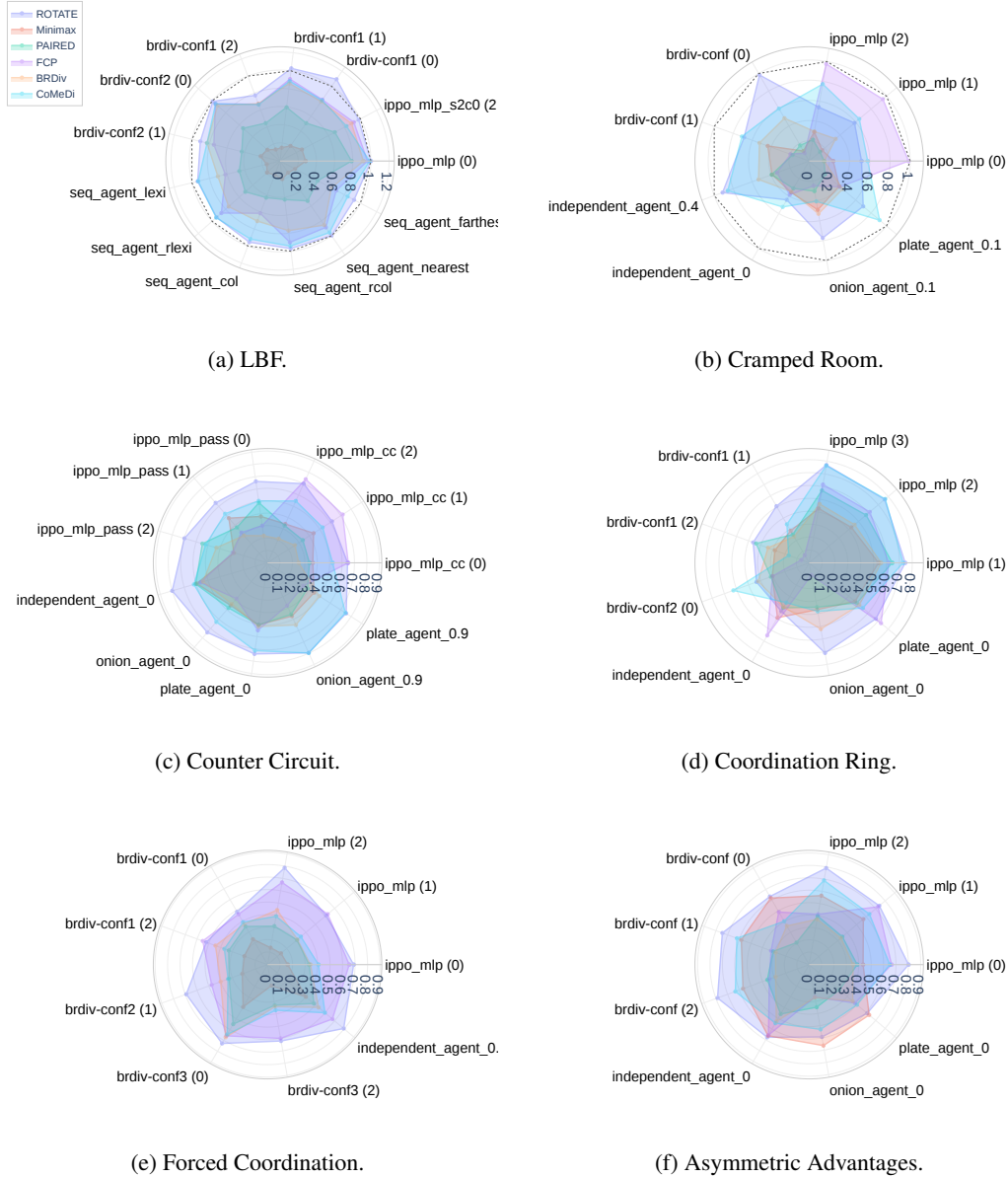


Figure 7: Normalized mean returns of ROTATE and all baselines across all tasks, broken down by evaluation teammate in Π^{eval} . Legend shown for LBF applies for all plots.

1032 number of foods, grid world size, level of observability, and whether to set the food level equal to the
 1033 sum to player levels in order to force players to coordinate to load each food.

1034 The experiments in this paper configured the LBF environment to a 7×7 grid, where two players
 1035 interact to collect three foods. Our LBF configuration is shown in Fig. 9. Each player observes the
 1036 full environment state, allowing each player to observe the locations of other agents and all foods
 1037 and the number of time steps elapsed in the current episode. Each player has six discrete actions: up,
 1038 down, left, right, no-op, and *load*, where the last action is the special food collection action. A food
 1039 may only be collected if the sum of player levels is greater than the level of the food. Since this paper
 1040 focuses on fully cooperative scenarios, we set the food level equal to the level of both players, so all
 1041 foods require cooperation in order to be collected. When a food is collected, both players receive an
 1042 identical reward, which is normalized such that the maximum return in an episode is 0.5. An episode

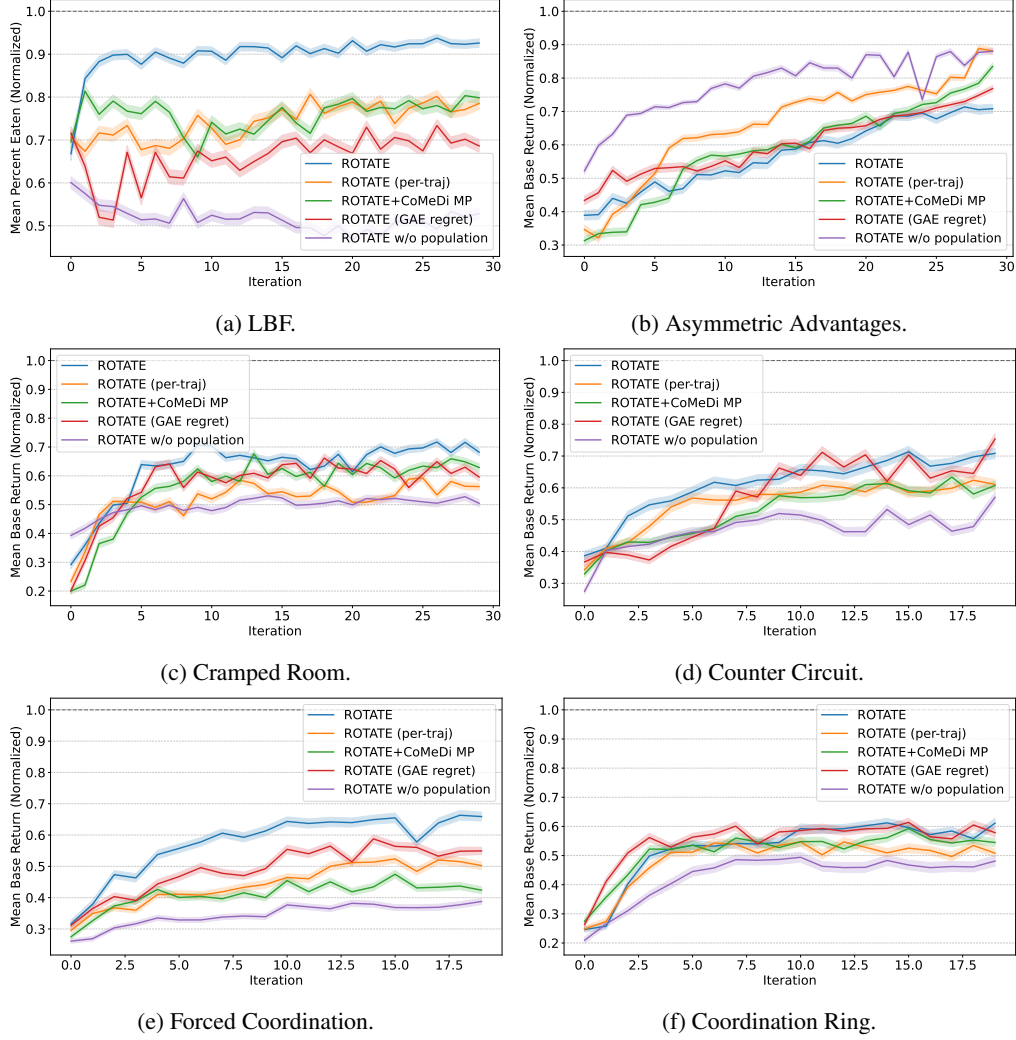


Figure 8: Learning curves of ROTATE and all variations of ROTATE considered in this paper. Normalized mean returns and bootstrapped 95% confidence intervals on Π^{eval} are shown.

1043 terminates if an invalid action is taken, players collide, or when 100 time steps have passed. Player
 1044 and food locations are randomized for each episode.

1045 **Overcooked** Introduced by Carroll et al. [9], the Overcooked suite
 1046 is a set of two-player collaborative cooking tasks, based on the
 1047 commercially successful Overcooked video game. Designed to study
 1048 human-AI collaboration, the original Overcooked suite consists of
 1049 five simple environment *layouts*, where two agents collaborate within
 1050 a grid world kitchen to cook and deliver onion soups. While Carroll
 1051 et al. [9] introduced Overcooked to study human-AI coordination,
 1052 Overcooked has become popularized for AHT research as well [10,
 1053 47, 19].

1054 We use the Jax re-implementation of the Overcooked suite by Ruther-
 1055 ford et al. [46], which is based on the original implementation by Car-
 1056 roll et al. [9]. Later versions of Overcooked include features such
 1057 as multiple dish types, order lists, and alternative layouts, but this
 1058 paper considers only the five original Overcooked layouts: Cramped

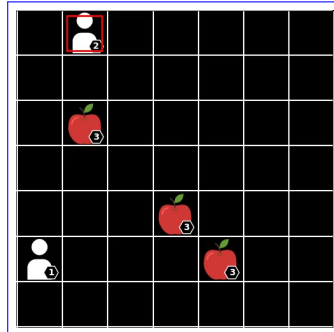


Figure 9: Level-based foraging environment. The apple icons denote food. The number on the icon indicates each player’s and food’s level. The AHT player is indicated by the red box.

1059 Room (CR), Asymmetric Advantages (AA), Counter Circuit (CC),
 1060 Coordination Ring (CoR), and Forced Coordination (FC).

1061 The objective for all five tasks is to deliver as many onion soups as
 1062 possible, where the only difference between the tasks is the environ-
 1063 ment layout, as shown in Fig. 10. To deliver an onion soup, players
 1064 must place three onions in a pot to cook, use a plate to pick up the
 1065 cooked soup, and send the plated soup to the delivery location. Each player observes the state and
 1066 location of all environment features (counters, pots, delivery, onions, and plates), the position and
 1067 orientation of both players, and an urgency indicator, which is 1 if there are 40 or fewer remaining
 1068 time steps, and 0 otherwise. Each player has six discrete actions, consisting of the four movement
 1069 actions, interact, and no-op. The reward function awards both agents +20 upon successfully deliver-
 1070 ing a dish, which is the return reported in the experimental results. To improve sample efficiency, all
 1071 algorithms are trained using a shaped reward function that provides each agent an additional reward
 1072 of 0.1 for picking up an onion, 0.5 for placing an onion in the pot, 0.1 for picking up a plate, and 1.0
 1073 for picking up a soup from the pot with a plate. An episode terminates after 400 time steps. Player
 1074 locations are randomized in each episode. In divided layouts such as AA and FC, we ensure that a
 1075 player is spawned on each half of the layout.

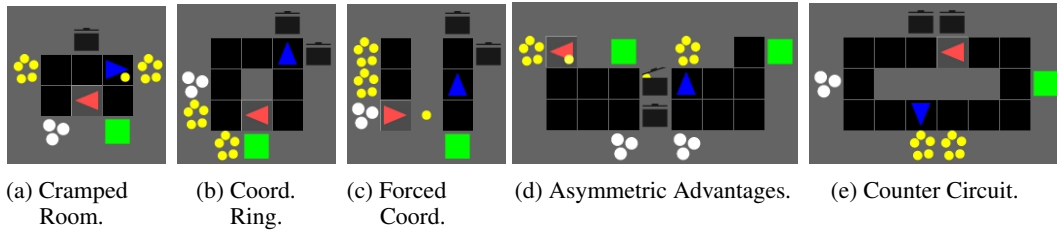


Figure 10: The five classic Overcooked layouts. Each yellow circle is an onion, while white circles are plates. Grid spaces with multiple yellow (resp. white) circles are onion (resp. plate) piles, which agents must visit to pick up an onion (or plate). The green square is the delivery location, where finished dishes must be sent to receive a reward. Black squares denote free space, while adjacent gray spaces are empty counters. A black pot icon indicates pots, while agents are shown as red and blue pointers. The AHT agent is highlighted.

1076 F Implementation Details

1077 As implementations of prior methods use PyTorch, but this project uses Jax, we re-implemented
 1078 all methods in this paper, using PPO [49] with Generalized Advantage Estimation (GAE) [48]
 1079 as a base RL algorithm and Adam [26] as the default optimizer. An anonymized version of the
 1080 code is released for reproducibility at <https://anonymous.4open.science/r/rotate/>, and we
 1081 ~~recommend consulting it for a full understanding of method implementations. Pseudocode for~~ [A](#)
 1082 ~~detailed discussion of ROTATE’s losses, update procedure, and pseudocode~~ is provided in App. A.
 1083 This section discusses implementation details such as training time choices, agent architectures, and
 1084 key hyperparameters for ROTATE and all baselines.

1085 F.1 Training Compute

1086 For fair comparison, all open-ended methods (ROTATE and all variations, PAIRED, Minimax
 1087 Return) were trained for the same number of open-ended learning iterations and a similar number of
 1088 environment interactions. For two-stage teammate generation approaches (FCP, BRDiv, CoMeDi), the
 1089 teammate generation stage is run using a similar amount of compute as the original implementations,
 1090 while the ego agent training stage is run for a sufficiently large number of steps to allow convergence.
 1091 We describe the amount of compute used for the teammate generation stage of each baseline below.

1092 In particular, the FCP population is generated by training 22-23 seeds of IPPO with 5 checkpoints
 1093 per seed for a population of approximately 110 agents—similar to Strouse et al. [52], who trained
 1094 32 seeds of IPPO with 3 checkpoints per seed for a population size of 96 agents. On the other hand,
 1095 BRDiv was trained with a population size of 3-4 agents, until we observed that each agent’s learning
 1096 converged. While we attempted training BRDiv with a larger population size, the algorithm was
 1097 prone to discovering degenerate solutions where only 2-3 agents in the population could discover

solutions with high SP returns, and all other agents in the population would have zero returns. Finally, CoMeDi was trained with a population size of 10 agents, until each agent’s learning converged. We attempted to train CoMeDi with a larger population size, but due to the algorithm’s quadratic complexity in the population size, its runtime surpassed the available time budget. Nevertheless, the population size of 10 forms a reasonable comparison to ROTATE because (1) the original paper used a population size of 8 for all Overcooked tasks, and (2) the configuration of CoMeDi in this paper runs for a similar wall-clock time as ROTATE.

F.2 Agent Architectures

For all methods considered in this paper, agents are implemented using neural networks and an actor-critic architecture, as is standard for PPO-based RL algorithms. All AHT methods implement policies without parameter sharing [13], to enable greater behavioral diversity. Specifics for ego agents, teammates, and best response agents are described below.

As mentioned in the main paper, ego agents are history-conditioned. Thus, ego agents are implemented with the S5 actor-critic architecture, a recently introduced recurrent architecture shown to have stronger long-term memory than prior types of recurrent architectures. Another advantage of the S5 architecture over typical recurrent architectures (e.g., LSTMs) is that it is parallelizable during training, allowing significant speedups in Jax [32].

On the other hand, teammates and best response agents are state-based. Best response agents are implemented with fully connected neural networks. Teammates are also based on fully connected neural networks, but the precise architecture varies based on the algorithm. For methods where the teammate only interacts with itself (FCP) or with the ego agent (Minimax Return), a standard actor-critic architecture is used. However, for open-ended learning methods that optimize regret (ROTATE and PAIRED), or for teammate generation methods that optimize adversarial diversity (ComeDi and BRDiv), teammates must estimate returns when interacting with multiple agents. Thus, for these methods, the teammate architecture includes a critic for each type of interaction.

In particular, for ROTATE and PAIRED, the teammate must estimate returns when interacting with the ego agent and its best response, and so it maintains a critic network for each partner type. For CoMeDi and BRDiv, given a population with n agents, each teammate must estimate the return when interacting with the other $n - 1$ agents in the population. As it would be impractical to maintain $n - 1$ critics for each teammate, the teammate instead uses a critic that conditions on the agent ID of a candidate partner agent—in effect, implementing the $n - 1$ critics via *parameter sharing* [13].

| Task | LBF | CR | AA | CC | CoR | FC |
|-------------|--------------------------------|--|--------------------|-------------------------|--------------------------|-------------------------|
| Timesteps | 3e5 , 1e6 | 1e6 | 1e6 | 1e6, 3e6 | 3e6 | 1e6, 3e6 , 1e7 |
| Number envs | 8 , 16 | 8 , 16 | 8 | 8 , 16 | 8 | 8 |
| Epochs | 7, 15 | 15 | 15 | 15 , 30 | 15 | 15 |
| Minibatches | 4 , 8 | 4, 8, 16 , 32 | 16 | 16 | 16 | 16 |
| Clip-Eps | 0.03 , 0.05 | 0.03, 0.05, 0.10, 0.15, 0.2 , 0.3 | 0.2, 0.3 | 0.1 , 0.2 | 0.1 , 0.2, 0.3 | 0.1 , 0.2 |
| Ent-Coef | 5e-3, 0.01 , 0.03, 0.05 | 5e-3, 0.01 , 0.03, 0.05 | 0.01 , 0.02 | 0.01, 0.03, 0.05 | 0.001, 0.01, 0.05 | 0.01, 0.05 |
| LR | 1e-4 | 1e-4 | 1e-4 , 1e-3 | 1e-4, 1e-3 | 1e-4, 5e-4, 1e-3 | 1e-4, 5e-4, 1e-3 |
| Anneal LR | true , false | true , false | true | true , false | true | true |

Table 3: Hyperparameters for IPPO.

| | LBF | CR | AA | CC | CoR | FC |
|-----------------|--------------------------|--------------------|--------------|-------------------|-------------------|-------------------------------|
| Timesteps | 4.5e7 | 4.5e7 | 4.5e7 | 9e7 | 9e7 | 9e7 |
| XP Coefficient | 0.1 , 0.75, 1, 10 | 1 , 10 | 10 | 0.01 , 10 | 0.01 , 10 | 0.01 , 0.1, 0.5, 1, 10 |
| Population size | 3 , 4, 5, 10 | 2, 3, 4 , 5 | 3 , 4 | 3 , 4 | 3 , 4 | 3 , 4 |
| Num Envs | 8, 32 | 8, 32 | 8, 32 | 8, 32 | 8, 32 | 8, 32 |
| LR | 1e-4, 5e-4 | 1e-4 | 1e-4 | 1e-3 | 1e-3, 5e-4 | 1e-3, 5e-4 |
| Ent-Coef | 0.01 | 0.01 | 0.01 | 0.05 | 0.05 | 0.05 |
| Clip-Eps | 0.03, 0.05 | 0.05 , 0.2 | 0.3 | 0.01 , 0.1 | 0.05, 0.1 | 0.05 , 0.1 |

Table 4: Hyperparameters for the teammate generation stage of BRDiv.

| | LBF | CR | AA | CC | CoR | FC |
|--------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------|--------------------------------|--------------------------------------|
| OEL Iterations | 30 | 30 | 30 | 20 | 20 | 20 |
| Num Envs | 16 | 16 | 16 | 16 | 16 | 16 |
| Regret-SP Weight | 1, 2 | 1, 3 | 1, 2 | 1, 2 | 1, 2 | 1, 2 |
| Minibatches | 4 , 8 | 8 | 8 | 8 | 8 | 8 |
| Timesteps per Iter (Ego) | 2e6 | 2e6 | 2e6 | 6e6 | 6e6 | 6e6 |
| Epochs (Ego) | 5, 10 , 20 | 10 , 15 | 10 | 10 | 10 | 5 , 10 |
| Ent-Coef (Ego) | 1e-4 , 1e-3, 0.01, 0.05 | 1e-4, 1e-3 , 1e-2 | 1e-3 , 0.01 | 1e-3 , 0.05 | 1e-3 , 0.05 | 1e-4 , 1e-3, 1e-2 |
| LR (Ego) | 5e-5 , 1e-4, 1e-3 | 1e-5, 3e-5, 5e-5 , 1e-4 | 1e-5, 3e-5, 5e-5 , 1e-4 | 3e-5, 5e-5 , 1e-3 | 1e-5, 3e-5 , 5e-5, 1e-3 | 8e-6, 1e-5 , 3e-5, 5e-5, 1e-4 |
| Eps-Clip (Ego) | 0.05, 0.1 | 0.1 , 0.2 | 0.1 , 0.3 | 0.1 | 0.1 | 0.1 |
| Anneal LR (Ego) | true, false | true, false | true, false | true, false | true, false | true, false |
| Timesteps per Iter (T) | 1e7 | 6e6 | 6e6 | 1.6e7 | 1.6e7 | 1.6e7 |
| Epochs (T) | 20 | 20 | 20 | 20 | 20 | 20 |
| Ent-Coef (T) | 0.05, 0.01 | 0.01 | 0.01 | 0.05 | 0.05 | 0.01, 0.05 |
| LR (T) | 1e-4 , 1e-3 | 1e-4 | 1e-4 | 1e-3 | 1e-3 | 1e-3 , 1e-4 |
| Clip-Eps (T) | 0.1 | 0.1, 0.2 | 0.3 | 0.1 | 0.1 | 0.1 , 0.2 |
| Anneal LR (T) | true , false | true , false | false | false | false | true, false |

Table 5: Hyperparameters for ROTATE. Hyperparameters specific to the teammate training process are marked by "(T)".

| | LBF | CR | AA | CC | CoR | FC |
|-----------------|-------|-------|------|------|------|------|
| Total Timesteps | 3e7 | 3e7 | 3e7 | 6e7 | 6e7 | 6e7 |
| Num Envs | 8 | 8 | 8 | 8 | 8 | 8 |
| LR | 5e-5 | 5e-5 | 5e-5 | 5e-5 | 3e-5 | 1e-5 |
| Epochs | 10 | 10 | 10 | 10 | 10 | 5 |
| Minibatches | 4 | 4 | 4 | 4 | 4 | 4 |
| Ent-Coef | 1e-4 | 1e-3 | 1e-3 | 1e-3 | 1e-3 | 1e-4 |
| Clip-Eps | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 |
| Anneal LR | false | false | true | true | true | true |

Table 6: Hyperparameters for PPO ego agent for all teammate generation methods.

| | LBF | CR | AA | CC | CoR | FC |
|---------------------|------|------|------|------|------|------|
| Timesteps Per Agent | 2e6 | 2e6 | 2e6 | 4e6 | 4e6 | 4e6 |
| Num Seeds | 23 | 23 | 23 | 22 | 22 | 22 |
| Num Checkpoints | 5 | 5 | 5 | 5 | 5 | 5 |
| Num Envs | 8 | 8 | 8 | 8 | 8 | 8 |
| LR | 1e-4 | 1e-4 | 1e-4 | 1e-3 | 1e-3 | 1e-3 |
| Epochs | 15 | 15 | 15 | 15 | 15 | 15 |
| Minibatches | 4 | 16 | 16 | 16 | 16 | 16 |
| Ent-Coef | 0.01 | 0.01 | 0.01 | 0.05 | 0.05 | 0.05 |
| Eps-Clip | 0.03 | 0.2 | 0.3 | 0.1 | 0.1 | 0.1 |
| Anneal LR | true | true | true | true | true | true |

Table 7: Hyperparameters for teammate generation stage of FCP.

| | LBF | CR | AA | CC | CoR | FC |
|-------------------------|-------|-------|-------|-------|-------|-------|
| Timesteps Per Iteration | 6e6 | 6e6 | 6e6 | 1e7 | 1e7 | 1e7 |
| Population Size | 10 | 10 | 10 | 10 | 10 | 10 |
| Num Envs | 16 | 16 | 16 | 16 | 16 | 16 |
| LR | 5e-4 | 1e-4 | 1e-4 | 1e-3 | 5e-4 | 5e-4 |
| Epochs | 15 | 15 | 15 | 15 | 15 | 15 |
| Minibatches | 8 | 8 | 8 | 8 | 8 | 8 |
| Ent-Coef | 1e-3 | 0.01 | 0.01 | 0.05 | 0.1 | 0.01 |
| Eps-Clip | 0.05 | 0.05 | 0.3 | 0.01 | 0.05 | 0.05 |
| Anneal LR | false | false | false | false | false | false |
| α | 0.2 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| β | 0.4 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 |

Table 8: Hyperparameters for the teammate generation stage of CoMeDi.

| | LBF | CR | AA | CC | CoR | FC |
|-----------------|-------|-------|-------|-------|-------|-------|
| Timesteps | 7.5e7 | 7.5e7 | 7.5e7 | 1.5e8 | 1.5e8 | 1.5e8 |
| Num Seeds | 5 | 5 | 5 | 5 | 5 | 5 |
| Num Checkpoints | 10 | 10 | 10 | 10 | 10 | 10 |
| Num Envs | 16 | 16 | 16 | 16 | 16 | 16 |
| LR | 1e-3 | 1e-4 | 1e-4 | 1e-3 | 1e-3 | 1e-3 |
| Epochs | 15 | 15 | 15 | 15 | 15 | 15 |
| Minibatches | 4 | 8 | 8 | 8 | 8 | 8 |
| Ent-Coef | 0.05 | 0.01 | 0.01 | 0.05 | 0.05 | 0.05 |
| Eps-Clip | 0.1 | 0.2 | 0.3 | 0.1 | 0.1 | 0.1 |
| Anneal LR | false | false | false | false | false | false |

Table 9: Hyperparameters for PAIRED.

| | LBF | CR | AA | CC | CoR | FC |
|--------------------------|-------|-------|-------|-------|-------|-------|
| OEL Iterations | 30 | 30 | 30 | 20 | 20 | 20 |
| Num Envs | 16 | 16 | 16 | 16 | 16 | 16 |
| Timesteps Per Iter (Ego) | 1e6 | 1e6 | 1e6 | 3e6 | 3e6 | 3e6 |
| Timesteps Per Iter (T) | 1e6 | 1e6 | 1e6 | 3e6 | 3e6 | 3e6 |
| LR | 1e-4 | 1e-4 | 1e-4 | 1e-3 | 1e-3 | 1e-3 |
| Epochs | 15 | 15 | 15 | 15 | 15 | 15 |
| Minibatches | 4 | 8 | 8 | 8 | 8 | 8 |
| Ent-Coef | 0.01 | 0.01 | 0.01 | 0.05 | 0.05 | 0.05 |
| Eps-Clip | 0.03 | 0.2 | 0.3 | 0.1 | 0.1 | 0.1 |
| Anneal LR | false | false | false | false | false | false |

Table 10: Hyperparameters for Minimax Return. Hyperparameters specific to the teammate training process are marked by "(T)".

1129 F.3 Hyperparameters

1130 This section presents the hyperparameters for ROTATE (Table 5), baseline methods (Tables 4 and 6
1131 to 10), and training evaluation teammates with IPPO (Table 3). Note that hyperparameters for the
1132 two-stage teammate generation methods are presented in separate tables, where those corresponding
1133 to the shared ego agent training stage are presented in Table 6. All experiments in the paper were
1134 performed with a discount factor of $\gamma = 0.99$ and $\lambda^{\text{GAE}} = 0.95$.

1135 Hyperparameters were searched for IPPO, BRDiv, and ROTATE, in that order, with the search for
1136 earlier methods informing initial hyperparameter values for later methods. Based on prior experience
1137 with PPO, we primarily searched the number of environments, epochs, minibatches, learning rate,
1138 entropy coefficient, the epsilon used for clipping the PPO objective, and whether to anneal the
1139 learning rate. For each hyperparameter, the searched values are listed in the tables, and selected
1140 values are bolded. We performed the search manually, typically varying one parameter over the listed
1141 range while holding others fixed, and varying parameters jointly only when varying one at a time did
1142 not yield desired results.

1143 Due to compute constraints, hyperparameters for FCP, CoMeDi, PAIRED, and Minimax Return were
1144 set based on knowledge of appropriate ranges gained from doing the hyperparameter searches over
1145 IPPO, BRDiv, and ROTATE.

| Name | Description | Est. BR Return |
|--------------------|--|----------------|
| brdiv_conf1(0) | Teammate trained by BRDiv. | 97.396 |
| brdiv_conf1(1) | - | 100.0 |
| brdiv_conf1(2) | - | 89.583 |
| brdiv_conf2(0) | - | 100.0 |
| brdiv_conf2(1) | - | 62.5 |
| ippo_mlp(0) | Teammate trained by IPPO to maximize return. | 100.0 |
| ippo_mlp_s2c0(2,0) | An intermediate checkpoint of a teammate trained by IPPO to maximize return. | 96.354 |
| seq_agent_col | Planning agent that collects food in column-major order (left to right, top to bottom). | 100.0 |
| seq_agent_rcol | Planning agent that collects food in reverse column-major order (right to left, bottom to top). | 100.0 |
| seq_agent_lexi | Planning agent that collects food in lexicographic order (top to bottom, left to right). | 100.0 |
| seq_agent_rlexi | Planning agent that collects food in reverse lexicographic order (bottom to top, right to left). | 100.0 |
| seq_agent_nearest | Planning agent that collects food in nearest to farthest order, based on the Manhattan distance from the agent's initial position. | 100.0 |
| seq_agent_farthest | Planning agent that collects food in farthest to nearest order, based on the Manhattan distance from the agent's initial position. | 100.0 |

Table 11: Evaluation teammates for LBF and estimated best response returns (percent eaten). Hyphens indicate that the agent description is the same as the previous description.

As described in Section 7 of the main paper, evaluation teammates were constructed using three strategies: training IPPO teammates in self-play using varied seeds and reward shaping, training teammates with BRDiv, and manually programming heuristic agents. Note that the evaluation teammates trained using IPPO and BRDiv were trained using different seeds than those used for training ROTATE and baseline methods.

The teammate construction procedure results in distinct teammate archetypes. Generally, IPPO agents execute straightforward, return-maximizing strategies. On the other hand, since BRDiv agents are trained to maximize self-play returns with their best response partner and to minimize cross-play returns with all other best response policies in the population, the generated teammates display more adversarial behavior compared to IPPO and heuristics. Coefficients on the SP and XP returns were carefully tuned to ensure that the behavior was not too adversarial, which we operationalized as teammates where the SP returns were high, but the XP returns were near zero.

Finally, the manually programmed heuristic agents have a large range of skills and levels of determinism. The LBF heuristics are planning-based agents that deterministically attempt to collect the apples in a specific order. Given a best response partner, the LBF heuristics can achieve the optimal task return in LBF. The Overcooked heuristics execute pre-programmed roles that are agnostic to the layout and some basic collision-avoidance logic. The "onion" heuristic collects onions and places them in non-full pots. The "plate" heuristic plates soups that are ready, and delivers them. The "independent" heuristic attempts to fulfill both roles by itself. All three heuristic types have a user-specified parameter that defines the probability that the agent places whatever it is holding on a nearby counter. The feature serves two purposes: first, it creates a larger space of behaviors, and second, it allows the heuristics to work for the FC task, where the agent in the left half of the kitchen must pass onions and plates to the right, while the agent in the right half must pick up resources from the dividing counter, cook soup, and deliver.

| Name | Description | Est. BR Return |
|-----------------------|---|----------------|
| brdiv_conf(0) | Teammate trained by BRDiv. | 214.063 |
| brdiv_conf(1) | - | 240.940 |
| ippo_mlp(0) | Teammate trained by IPPO to maximize return. | 256.875 |
| ippo_mlp(1) | - | 253.750 |
| ippo_mlp(2) | - | 249.686 |
| independent_agent_0.4 | Agent programmed to cook and deliver soups. If holding item, 40% chance of placing item on the counter. | 197.188 |
| independent_agent_0 | Agent programmed to cook and deliver soups. | 132.50 |
| onion_agent_0.1 | Agent programmed to place onions in non-full pots. If holding item, 10% chance of placing item on counter. | 146.875 |
| plate_agent_0.1 | Agent programmed to plate finished soups and deliver. If holding item, 10% chance of placing item on counter. | 191.250 |

Table 12: Evaluation teammates for Cramped Room and estimated best response returns. Hyphens indicate that the agent description is the same as the previous description.

| Name | Description | Est. BR Return |
|---------------------|--|----------------|
| brdiv_conf(0) | Teammate trained by BRDiv. | 286.875 |
| brdiv_conf(1) | - | 335.625 |
| brdiv_conf(2) | - | 333.750 |
| ippo_mlp(0) | Teammate trained by IPPO to maximize return. | 382.50 |
| ippo_mlp(1) | - | 369.375 |
| ippo_mlp(2) | - | 312.50 |
| independent_agent_0 | Agent programmed to cook and deliver soups. | 308.125 |
| onion_agent_0 | Agent programmed to place onions in non-full pots. | 301.250 |
| plate_agent_0 | Agent programmed to place onions in non-full pots. | 285.0 |

Table 13: Evaluation teammates for Asymmetric Advantages and estimated best response returns. Hyphens indicate that the agent description is the same as the previous description.

1171 Descriptions of the evaluation teammates for each task and estimated best response returns are
1172 provided in Tables 11 to 16.

1173 **Evaluation Return Normalization Details.** The lower return bound is set to zero since a poor
1174 teammate could always cause a zero return in all tasks considered. Ideally, the upper return bounds
1175 would be the returns achieved with the theoretically optimal best response teammate for each
1176 evaluation teammate. To approximate this, we instead set the upper bound equal to the maximum
1177 average return achieved by any method, for each evaluation teammate.

1178 As described in Section 7, our normalized return metric is similar to the BRProx metric recommended
1179 by Wang et al. [58]. The main difference is that we aggregate results using the mean rather than
1180 the interquartile mean (IQM), due to challenges around determining appropriate upper bounds for
1181 return normalization. In particular, during method development, we used looser BR return estimates
1182 to perform return normalization, leading to normalized returns often surpassing 1.0 for certain
1183 teammates. Under such conditions, aggregating results using the IQM led to entirely dropping results
1184 corresponding to particular teammates.

| Name | Description | Est. BR Return |
|---------------------|---|----------------|
| ippo_mlp_cc(0) | Teammate trained by IPPO to maximize return. Navigates counterclockwise around counter. | 200.625 |
| ippo_mlp_cc(1) | - | 198.120 |
| ippo_mlp_cc(2) | - | 194.375 |
| ippo_mlp_pass(0) | Teammate trained by IPPO+reward shaping to pass onions across the counter. | 137.813 |
| ippo_mlp_pass(1) | - | 103.125 |
| ippo_mlp_pass(2) | - | 170.0 |
| independent_agent_0 | Agent programmed to cook and deliver soups. | 77.189 |
| onion_agent_0.9 | Agent programmed to place onions in non-full pots. If holding item, 90% chance of placing item on counter. | 80.0 |
| onion_agent_0 | Agent programmed to place onions in non-full pots. | 81.563 |
| plate_agent_0.9 | Agent programmed to plate finished soups and deliver. If holding item, 90% chance of placing item on counter. | 97.189 |
| plate_agent_0 | Agent programmed to place onions in non-full pots. | 76.875 |

Table 14: Evaluation teammates for Counter Circuit and estimated best response returns. Hyphens indicate that the agent description is the same as the previous description.

| Name | Description | Est. BR Return |
|---------------------|--|----------------|
| brdiv_conf1(1) | Teammate trained by BRDiv. | 161.250 |
| brdiv_conf1(2) | - | 183.440 |
| brdiv_conf2(0) | - | 142.810 |
| ippo_mlp(1) | Teammate trained by IPPO to maximize return. | 249.688 |
| ippo_mlp(2) | - | 246.560 |
| ippo_mlp(3) | - | 246.560 |
| independent_agent_0 | Agent programmed to cook and deliver soups. | 136.250 |
| onion_agent_0 | Agent programmed to place onions in non-full pots. | 72.50 |
| plate_agent_0 | Agent programmed to place onions in non-full pots. | 110.938 |

Table 15: Evaluation teammates for Coordination Ring and estimated best response returns. Hyphens indicate that the agent description is the same as the previous description.

| Name | Description | Est. BR Return |
|-----------------------|---|----------------|
| brdiv_conf1(0) | Teammate trained by BRDiv. | 131.560 |
| brdiv_conf1(2) | - | 184.690 |
| brdiv_conf2(1) | - | 143.750 |
| brdiv_conf3(0) | - | 71.250 |
| brdiv_conf3(2) | - | 174.690 |
| ippo_mlp(0) | Teammate trained by IPPO to maximize return. | 220.0 |
| ippo_mlp(1) | - | 214.380 |
| ippo_mlp(2) | - | 225.620 |
| independent_agent_0.6 | Agent programmed to cook and deliver soups. If holding item, 60% chance of placing item on the counter. | 81.250 |

Table 16: Evaluation teammates for Forced Coordination and estimated best response returns. Hyphens indicate that the agent description is the same as the previous description.

H Compute infrastructure

Experiments were performed on two servers, each with the following specifications:

- **CPUs:** two Intel(R) Xeon(R) Gold 6342 CPUs, each with 24 cores and two threads per core.
- **GPUs:** four NVIDIA A100 GPUs, each with 81920 MiB VRAM.

The experiments in this paper were implemented in Jax and parallelized across seeds. On the servers above, each method took approximately 4-6 hours of wall-clock time to run.