
A Minimax-Bayes Approach to Ad Hoc Teamwork

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning policies for Ad Hoc Teamwork (AHT) is challenging. Most standard
2 methods choose a specific distribution over training partners, which is assumed to
3 mirror the distribution over partners after deployment. Moreover, they offer limited
4 guarantees over worst-case performance. To tackle the issue, we propose using a
5 worst-case prior distribution by adapting ideas from minimax-Bayes analysis to
6 AHT. We thereby explicitly account for our uncertainty about the partners at test
7 time. Extensive experiments, including evaluations on coordination tasks from the
8 Melting Pot suite, show our method’s superior robustness compared to self-play,
9 fictitious play, and best response learning w.r.t. policy populations. This highlights
10 the importance of selecting an appropriate training distribution over teammates to
11 achieve robustness in AHT.

12 1 Introduction

13 Domain generalisation is often crucial in Reinforcement Learning (RL) and is typically assessed by
14 placing an agent in novel environments (Cobbe et al., 2019). Likewise, in Multi-Agent Reinforcement
15 Learning (MARL), generalisation to new agents can be evaluated by pairing a trained policy with
16 unseen actors (Barrett et al., 2011; Hu et al., 2020; Leibo et al., 2021; Agapiou et al., 2023). While
17 zero-shot domain adaptation is a valuable property (Higgins et al., 2017; Schäfer, 2022), it is equally
18 important to ensure proper transfer to new behaviours in multi-agent settings, especially in situations
19 where undesired interactions may arise (Gleave et al., 2019). More specifically, Ad Hoc Teamwork
20 (AHT) occurs when multiple agents, initially unfamiliar with each other, must collaborate to achieve
21 a common goal. In a world where autonomous agents are being progressively introduced in such
22 tasks, cooperation with humans is becoming a major concern (Stone et al., 2010; Ji et al., 2023).

23 Efforts in AHT have primarily focused on learning and inferring models of teammates’ behaviours
24 (Barrett et al., 2011; Albrecht et al., 2015; Barrett et al., 2017; Chen et al., 2020; Muglich et al., 2022b),
25 adapting to behaviour shifts (Ravula et al., 2019), and enhancing generalisation by encouraging
26 diversity in partners during training (Jaderberg et al., 2019; Hu et al., 2020; Charakorn et al., 2020;
27 Lupu et al., 2021; Strouse et al., 2021). However, these methods provide limited guarantees regarding
28 worst-case AHT performance.

29 A multi-agent system can encompass numerous and diverse *scenarios*, each characterised by its
30 actors. For example, autonomous cars operate alongside various human drivers and other autonomous
31 vehicles. Similarly, in a surgical setting, a robot may need to cooperate with surgeons who have a
32 wide range of habits and expertise levels. In each of these scenarios, we can adopt the perspective that
33 the *focal* actors are controlled by the learner, whereas the other actors are viewed as fixed, forming
34 the *background* of the task (Leibo et al., 2021; Agapiou et al., 2023). These scenarios can be viewed
35 as distinct single-agent environments, as each combination of background actors induces different
36 transition dynamics and reward functions. A common practice involves constructing representative
37 scenarios and training a policy on a uniform distribution over them (Strouse et al., 2021; Lupu et al.,
38 2021). However, this only ensures good performance for that specific distribution.

39 Recent studies in zero-shot domain transfer showed that selecting an appropriate prior over training
40 environments is key to learning robust policies (Pinto et al., 2017; Dennis et al., 2020; Garcin et al.,
41 2023; Jiang et al., 2021; Buening et al., 2023). Intuitively, this insight should apply to the AHT
42 setting as well, suggesting that choosing a specific prior over scenarios/partners may improve the
43 robustness of learned policies. Assuming that no information is available about the teammates at test
44 time (and their distribution), we consider the *worst* possible prior over the set of partners given our
45 policy, an idea adopted from the minimax-Bayes concept (Berger, 1985).

46 **Contributions.** We summarise our contributions as follows:

- 47 1. We adapt Minimax-Bayes Reinforcement Learning (MBRL)(Buening et al., 2023) to the AHT
48 setting, reasoning about uncertainty with respect to partners rather than environments (Section 4).
- 49 2. We examine the advantages of using utility and regret to measure performance in the AHT setting,
50 and propose algorithms to target either metric (Section 5).
- 51 3. We adapt a Gradient Descent-Ascent (GDA) (Lin et al., 2020) based algorithm, in conjunction with
52 policy-gradient methods, and discuss its convergence guarantees for softmax policies (Section 6).
- 53 4. We conduct extensive experiments to evaluate our approach. We test learned policies on both seen
54 and held-out scenarios for various cooperative problems, including partially observable games
55 such as environments from the Melting Pot suite (Leibo et al., 2021; Agapiou et al., 2023). We
56 compare our approach against Self-Play (SP), Fictitious Play (FP)(Brown, 1951; Heinrich et al.,
57 2015) as well as learning a policy on a fixed uniform distribution over scenarios (Lupu et al.,
58 2021), which is closely related to Fictitious Co-Play (FCP) (Strouse et al., 2021), as both learn the
59 best response to population of policies (Section 7).
- 60 5. Our results confirm the theory and empirically demonstrate that our approach leads to the most
61 robust solutions for both simple and deep RL coordination tasks, even when teammates are
62 adaptive. This highlights the importance of choosing an appropriate distribution over training
63 scenarios to develop policies that better transfer to new teammates.

64 2 Related Work

65 **Ad Hoc Teamwork.** In AHT, we are interested in developing agents capable of cooperating with
66 other unfamiliar agents without any form of prior coordination (Rovatsos and Wolf, 2002; Stone et al.,
67 2010; Barrett et al., 2011, 2017). Popular approaches usually involve some form of Population Play
68 (PP), where policies forming a population are learning by interacting with each other (Lupu et al.,
69 2021; Muglich et al., 2022a; Leibo et al., 2021; Agapiou et al., 2023). Key strategies for ensuring
70 generalisation to new partners include promoting policy diversity within the training population
71 (Charakorn et al., 2020) and preventing overfitting to training partners (Lanctot et al., 2017). Both
72 Lupu et al. (2021) and Strouse et al. (2021) previously showed that learning a best response to a more
73 diverse population leads to improved generalisation. Additionally, Jaderberg et al. (2019) showed the
74 effectiveness of PP when diversity is encouraged through evolving pseudo-rewards. However, PP still
75 struggles with producing policies that are robustly collaborative with new partners and sometimes
76 exhibits overfitting (Carroll et al., 2019; Leibo et al., 2021; Agapiou et al., 2023).

77 To push the boundaries of AHT further, many studies use inference on teammate models to maintain
78 a belief about ad hoc partners based on previous interactions within an episode (Barrett et al., 2011;
79 Albrecht et al., 2015). Efforts have also been made to improve the learning and generalisation of such
80 models to new partners (Barrett and Stone, 2015; Barrett et al., 2017; Muglich et al., 2022b).

81 An alternative approach proposed by Li et al. (2019) involves a robust formulation of deep determin-
82 istic policy gradients, assuming worst-case teammates. Unlike our setup, they train a joint policy that
83 remains consistent throughout learning, and design their algorithm specifically for deep deterministic
84 policy gradients, while our approach is compatible with any policy-gradient algorithm.

85 Even though the aforementioned methods attempt at improving cooperative robustness, they always
86 assume specific distributions for the partners. Jaderberg et al. (2019) used a distribution favoring the
87 matchmaking of policies of similar levels under the intuition that the reward signal is stronger in
88 those cases, it does not provide any insights on its eventual effects on AHT robustness. As such, the
89 actual impact of training partner distribution on robustness is left under-explored and represents a
90 component that can be further exploited in conjunction with other AHT mechanisms.

91 **Zero-shot Domain Transfer.** Robustness to unknown partners can be seen as a form of zero-shot
92 domain transfer. Each possible team composition involving the agent of interest can be considered
93 a different environment. In the single agent setting, Jiang et al. (2021) demonstrated that adapting
94 the training environment distribution by prioritising environments with higher prediction loss (a
95 measure of the policy’s lack of knowledge) leads to improved sample efficiency and generalisation.
96 Building on this idea, Garcin et al. (2023) prioritised environments where the mutual information
97 between the learning policy’s internal representation and the environment identity was lower, using
98 information theory to achieve similar results. The idea of tempering with the environment distribution
99 was also explored by Pinto et al. (2017), who employed a maximin utility formulation to choose
100 continuous adversarial environment perturbations throughout learning. Instead of utility, Dennis
101 et al. (2020) stressed the advantages of using regret by proposing a training environment sampling
102 scheme avoiding entirely unsolvable and uninformative environments. Most interestingly, Buening
103 et al. (2023) conducted a study over worst-case priors (for both utility and regret) over training
104 environments, and proved that worst-case distributions are a good fit for domain transfer. Finally,
105 there exist works on domain transfer in the MARL setting (Schäfer, 2022), but this differs from
106 our focus on transferring to new partners. This related work is consistently in favor of caring about
107 environment distributions for robustness, providing strong motivation to bring this concern to AHT.

108 3 Problem Formulation

109 3.1 Preliminaries

110 An m -player Partially Observable Markov Game (POMG) is given by a tuple $\mu =$
111 $\langle \mathcal{S}, \mathcal{X}, \mathcal{A}, \mathcal{O}, P, \rho, T \rangle$ defined on finite sets of states \mathcal{S} , observations \mathcal{X} and actions \mathcal{A} . The ob-
112 servation function $\mathcal{O} : \mathcal{S} \times \{1, \dots, m\} \rightarrow \mathcal{X}$ provides a state space view for each player. In each
113 state, each player i chooses an action $a_i \in \mathcal{A}$. Following their joint action $\mathbf{a} = (a_1, \dots, a_m) \in \mathcal{A}^m$,
114 the state is updated according to the transition function $P : \mathcal{S} \times \mathcal{A}^m \rightarrow \Delta(\mathcal{S})$. After a transition,
115 each player receives a reward defined by $\rho : \mathcal{S} \times \mathcal{A}^m \times \{1, \dots, m\} \rightarrow \mathbb{R}$. The game ends after T
116 transitions. Permuting player indices does not have any effect on μ .

117 A policy $\pi : \mathcal{X} \times \mathcal{A} \times \mathcal{X} \times \mathcal{A} \times \dots \times \mathcal{X} \rightarrow \Delta(\mathcal{S})$ is a probability distribution over a single agent’s
118 actions, conditioned on that agent’s history of observations and actions. We denote Π the set of all
119 policies and $\Pi^D \subset \Pi$ the set of deterministic policies.

120 3.2 Scenarios

121 Let a *scenario* σ_b^c be defined by its number of *focal* players c , and its *background* players $\pi^b =$
122 $(\pi_1^b, \dots, \pi_{m-c}^b) \in \Pi^{m-c}$. We say we deploy a policy π^f in scenario σ_b^c if the c focal players are set
123 to copies of π^f . Hence, in addition to the $m - c$ many background policies π^b , there are c many
124 focal policies $\pi^f = (\pi^f, \dots, \pi^f)$. We sometimes use σ as a shorthand notation for σ_b^c for simplicity.
125 We also denote $\mathbf{a}^f \in \mathcal{A}^c$ and $\mathbf{a}^b \in \mathcal{A}^{m-c}$ the joint actions of the focal and background players,
126 respectively. A background population $\mathcal{B} \subset \Pi$ is a finite set of policies, to which we assign a set of
127 scenarios:

$$\Sigma(\mathcal{B}) \triangleq \{\sigma_b^c \mid 1 < c \leq m, \pi^b \in \mathcal{B}^{m-c}\}.$$

128 A scenario σ_b^c on μ can be viewed as its own c -player POMG, through the marginalisation of the
129 policies of its background players.¹ We denote $\mu(\sigma) = \langle \mathcal{S}, \mathcal{X}, \mathcal{A}, \mathcal{O}_\sigma, P_\sigma, \rho_\sigma, \gamma, T \rangle$ the POMG
130 induced by scenario σ , where $\mathcal{O}_\sigma : \mathcal{S} \times \{1, \dots, c\} \rightarrow \mathcal{X}$ is the corresponding observation function,
131 $P_\sigma : \mathcal{S} \times \mathcal{A}^c \rightarrow \Delta(\mathcal{S})$ the transition function given by

$$P_\sigma(s' \mid s, \mathbf{a}^f) = \begin{cases} P(s' \mid s, \mathbf{a}^f), & c = m \\ \sum_{\mathbf{a}^b} \left(P(s' \mid s, \mathbf{a}^f, \mathbf{a}^b) \prod_i \pi_i^b(\mathbf{a}_i^b \mid h_i) \right), & c < m \end{cases}$$

132 and $\rho_\sigma : \mathcal{S} \times \mathcal{A}^c \times \{1, \dots, c\} \rightarrow \mathbb{R}$ the induced reward function with:

$$\rho_\sigma(s, \mathbf{a}^f, i) = \begin{cases} \rho(s, \mathbf{a}^f, i), & c = m \\ \sum_{\mathbf{a}^b} \left(\rho(s, \mathbf{a}^f, \mathbf{a}^b, i) \prod_i \pi_i^b(\mathbf{a}_i^b \mid h_i) \right), & c < m \end{cases}$$

¹Each scenario can be seen as a decentralized partially observable Markov decision process (Oliehoek, 2012) constrained by the fact that the c players are copies.

133 where h_i is the history of observations and actions of the i -th policy and \mathbf{a}_i^b its action in \mathbf{a}^b . We denote
 134 the scenario that only involves copies of the focal policy, i.e. the universalisation scenario (Leibo
 135 et al., 2021), by $\sigma^{\text{SP}} = \sigma_{\emptyset}^m$.

136 3.3 Evaluation

137 The expected utility of a policy in scenario σ is the mean return of the focal policies given by the
 138 expected *focal-per-capita return* (Leibo et al., 2021; Agapiou et al., 2023):

$$U(\pi, \sigma) \triangleq \sum_{t=1}^T \frac{1}{c} \sum_{i=1}^c \mathbb{E}_{\mu(\sigma)}^{\pi} [\rho_{\sigma}(s_t, \mathbf{a}_t^f, i)]. \quad (1)$$

139 $U^*(\sigma) \triangleq \max_{\pi \in \Pi} U(\pi, \sigma)$ denotes the maximal utility achievable in scenario σ . This definition for
 140 utility represents the need for autonomous agents to always maximise the mean joint rewards of its
 141 copies, regardless of the scenario. We can further define the notion of regret incurred by deploying
 142 some policy π on scenario σ , as the gap between the maximal utility and the utility of π on σ :

$$R(\pi, \sigma) \triangleq U^*(\sigma) - U(\pi, \sigma). \quad (2)$$

143 To assess a learning method in terms of AHT, we use the evaluation protocol of Leibo et al. (2021).
 144 This has two phases:

- 145 1. **Training phase:** A test background population $\mathcal{B}^{\text{test}}$ is kept hidden. The policy learner has access
 146 to the game μ with no restriction, beside accessing $\mathcal{B}^{\text{test}}$. For example, the learner is free to use a
 147 modified instance μ' of μ , where \mathcal{O} could be changed to include observations of other players, or
 148 again where ρ could be tweaked to return the joint rewards rather than individual rewards.
- 149 2. **Testing phase:** The obtained policy is frozen and cannot be trained any further. We compute the
 150 performance of the policy on μ by taking its average expected utility across a series of held-out
 151 test scenarios $\Sigma^{\text{test}} \subset \Sigma(\mathcal{B}^{\text{test}})$. In addition to performance, we consider two metrics related to
 152 robustness, worst-case utility and worst-case regret:

$$p(\pi, \Sigma) = \frac{1}{|\Sigma|} \sum_{\sigma \in \Sigma} U(\pi, \sigma), \quad U^-(\pi, \Sigma) = \min_{\sigma \in \Sigma} U(\pi, \sigma), \quad R^+(\pi, \Sigma) = \max_{\sigma \in \Sigma} R(\pi, \sigma). \quad (3)$$

153 Maximising U^- is typically preferable when falling below a certain utility threshold must be
 154 avoided at all costs; for instance minimising casualties in a surgical context. Conversely, minimising
 155 R^+ avoids decisions that lead to significantly worse outcomes than the best-case.

156 The end objective is to design a learning process outputting a policy that reliably maximises its
 157 expected utility (focal-per-capita return) on possibly unseen scenarios.

158 3.4 Assumptions

159 To ensure our setting aligns with the AHT literature, we must adhere to three assumptions (Mirsky
 160 et al., 2022): a) the absence of prior coordination. The learner must be capable of cooperating with
 161 the team on-the-fly, without relying on previously established collaboration strategies, even between
 162 copies of the learner’s agent. b) There is no control over teammates, the learner can control its
 163 own copies but not other agents in the configuration. c) All agents are assumed to share a common
 164 objective. Nonetheless, their reward function may be different, reflecting varying preferences. In
 165 this work, we choose to address this last point by assuming a class of possible reward functions for
 166 the background players. In an attempt to model realistic situations, we formalise this diversity by
 167 considering various levels of prosociality λ (Peysakhovich and Lerer, 2017) and risk-aversion δ for
 168 each policy. To illustrate with an example, in a setting where company coworkers have to realise a
 169 project, there might be workers that have a high preference over their own contribution (with a better
 170 chance to get promoted later), while there may be others that are inclined to delegate their work for
 171 things they are unsure about to the team:

$$\rho_{\text{social+risk}}(s, \mathbf{a}, i) = \rho_{\text{social}}^+(s, \mathbf{a}, i) - \delta_i \rho_{\text{social}}^-(s, \mathbf{a}, i),$$

172 with ρ_{social} defined as

$$\rho_{\text{social}}(s, \mathbf{a}, i) = \lambda_i \rho(s, \mathbf{a}, i) + (1 - \lambda_i) \sum_{j=1}^m \rho(s, \mathbf{a}, j),$$

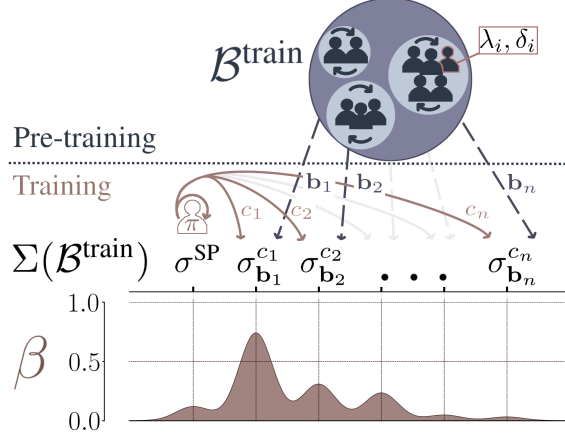


Figure 1: Comprehensive illustration of the framework used in this paper. Prior to training the focal policy π , background policies with different preferences (λ_i, δ_i) learn by interacting within sub-populations of varying sizes. These sub-populations are then combined to form a background population, $\mathcal{B}^{\text{train}}$, used as a common ‘train dataset’ for all algorithms.

Our primary focus is on the training phase, where the main policy π is learned alongside the distribution β over scenarios. These scenarios mix copies of π with policies from $\mathcal{B}^{\text{train}}$, where the self-play scenario σ^{SP} has the policy interacting only with copies of itself.

173 where f^+ and f^- are the positive and negative parts of f , and (λ_i, δ_i) are the levels of prosociality
 174 and risk-aversion for agent i . Combining values of prosociality and risk-aversion allows for the
 175 consideration of behaviours with a wide range of preferences.

176 4 Achieving Robust AHT

177 To learn a policy able to cooperate with new partners, a straightforward idea is to reconstruct scenarios
 178 that would likely be encountered in nature. A roadblock to this approach however is that it requires
 179 two main ingredients: a) a diverse pool of partners, and b) a prior distribution over them. The prior,
 180 often neglected, is important as it captures our uncertainty about the true partners observed in nature.

181 In Section 4.1, we reflect on motivating previous work on diverse behaviour generation, before
 182 describing our own adopted approach. Section 4.2 then introduces the Minimax-Bayes idea to AHT,
 183 by stating the existing connections of the setting with MBRL’s.

184 4.1 Constructing Training Scenarios

185 Prior to learning any robust policy, we need to construct diverse scenarios. A background population
 186 that encompasses a wide range of behaviours is needed. Previous work on AHT tackled the issue
 187 in various manners, such as using genetic algorithms (Muglich et al., 2022b), rule-based policies
 188 generated with MAP Elites (Canaan et al., 2023), SP policies (Strouse et al., 2021), explicit behavior
 189 diversification through regularisation (Lupu et al., 2021), or through evolved pseudo-rewards (Jader-
 190 berg et al., 2019). Based on real-life examples and aiming to thoroughly assess the effects of partner
 191 priors, we adopt the following approach:

- 192 • Each background policy has unique preferences (λ_i, δ_i) .
- 193 • Policies are organized into sub-populations $\mathcal{B} = \bigcup_k \mathcal{B}_k$ of varying sizes, simulating different
 194 communities.
- 195 • Each sub-populations are separately trained using PP. Given the diverse preferences and varying
 196 sizes of these sub-populations, distinct habits, common practices, and established conventions will
 197 emerge within each group, effectively mimicking various cultures.

198 This choice for constructing scenarios is rather arbitrary and is not the main focus of our work.
 199 Nevertheless, a rigorous generation procedure is important to bring forward the effects of various
 200 scenario priors on AHT robustness.

201 4.2 Minimax-Bayes AHT

202 In the standard single-agent Bayesian RL setting, the learner selects a subjective belief β over
 203 candidate Markov Decision Processes (MDPs) \mathcal{M} for the unknown, true environment $\mu^* \in \mathcal{M}$.
 204 The learner’s objective is to maximise its expected utility with respect to the chosen prior
 205 $U(\pi, \beta) = \int_{\mathcal{M}} \bar{U}(\pi, \mu) d\beta(\mu)$, i.e. finding the Bayes-optimal policy. In MBRL, Buening et al. (2023)
 206 proposed considering the worst possible prior for the agent, without knowledge of the policy that
 207 will be chosen. This approach can be interpreted as nature playing the minimising player against the
 208 policy learner in a simultaneous-move zero-sum normal-form game. Learning against a worst-case
 209 prior intuitively makes the learner more robust, as it prepares for the worst outcomes.

210 To transfer this idea to our setting, we remark that any finite population \mathcal{B} provides a finite set of
 211 POMGs $\mathcal{M}_{\mathcal{B}} = \{\mu(\sigma) | \sigma \in \Sigma(\mathcal{B})\}$. The difference here is the use of POMGs rather than MDPs. We
 212 extend the notion of expected utility with respect to a prior over scenarios, i.e. when $\beta \in \Delta(\Sigma(\mathcal{B}))$:

$$U(\pi, \beta) \triangleq \mathbb{E}_{\sigma \sim \beta}[U(\pi, \sigma)] = \sum_{\sigma} U(\pi, \sigma) \beta(\sigma).$$

213 This allows us to formulate the following maximin game:

$$\max_{\pi \in \Pi} \min_{\beta \in \Delta(\Sigma(\mathcal{B}))} U(\pi, \beta). \quad (4)$$

214 Similarly to Buening et al. (2023), we are interested in knowing whether such a game has a solution
 215 (i.e., a value), assuming that nature and the agent play simultaneously without knowledge of each
 216 other’s move. This is relevant in our setting because the policy learner does not know the true
 217 distribution of partners available in nature, while the effective nature’s distribution of scenarios should
 218 not depend on the agent’s policy. Fortunately, (4) has a value when \mathcal{B} is finite.

219 **Corollary 1.** *For an m -player POMG μ in a finite state-action space, with a known reward function
 220 and a finite horizon, and a finite background population \mathcal{B} , the maximin game (4) has a value:*

$$\max_{\pi \in \Pi} \min_{\beta \in \Delta(\Sigma(\mathcal{B}))} U(\pi, \beta) = \min_{\beta \in \Delta(\Sigma(\mathcal{B}))} \max_{\pi \in \Pi} U(\pi, \beta). \quad (5)$$

221 *Proof.* First, observe that for any stochastic policy $\pi \in \Pi$, there exists a distribution over deterministic
 222 policies $\phi \in \Delta(\Pi^{\text{D}})$ such that $\pi(a_t | h_t) = \sum_{d \in \Pi^{\text{D}}} d(a_t | h_t) \phi(d)$. Consequently, we can rewrite the
 223 utility as $U(\pi, \beta) = \sum_{d \in \Pi^{\text{D}}} \sum_{\sigma \in \Sigma(\mathcal{B})} U(d, \sigma) \phi(d) \beta(\sigma)$. This demonstrates that U is bilinear in ϕ
 224 and β , which allows us to apply the minimax theorem, thus proving the result. \square

225 Importantly, prior work that chooses an arbitrarily fixed prior is limited in terms of robustness
 226 guarantees: it only ensures maximal utility for their specific prior. In contrast, a policy π_U^* solving
 227 the maximin utility problem (4) has its expected utility lower-bounded on $\Sigma(\mathcal{B})$:

$$\forall \beta \in \Delta(\Sigma(\mathcal{B})), \quad U(\pi_U^*, \beta) \geq U(\pi_U^*, \beta_U^*), \quad (6)$$

228 where β_U^* is the worst-case prior for π_U^* . Simply put, π_U^* performs the worst when the prior is its
 229 worst-case β_U^* , but can only improve when the prior deviates from β_U^* . Additionally, it is also optimal
 230 on the worst-case prior:

$$\forall \pi \in \Pi, \quad U(\pi_U^*, \beta_U^*) \geq U(\pi, \beta_U^*). \quad (7)$$

231 Note that is is entirely different from the best response to the fixed worst-case prior
 232 $\arg \max_{\pi} U(\pi, \beta_U^*)$, which once again, only has a guaranteed optimal utility on β_U^* .

233 5 Utility or Regret?

234 Optimising for the worst-case utility (4) might be problematic. Nature could resort to only picking
 235 scenarios where the the focal players achieve the worst possible score. Then, the prior trivially
 236 minimises utility for any chosen policy. Buening et al. (2023) addresses this issue by instead
 237 considering the regret of a policy.. The difference is that ‘impossible’ scenarios will always yield zero
 238 regret for any policy, thus becoming irrelevant for a regret-maximising nature. Letting $L(\pi, \beta) \triangleq$
 239 $\int_{\mathcal{M}} R(\pi, \mu) d\beta(\mu)$ be the Bayesian regret with respect to a prior β , we now formulate the following
 240 minimax regret game:

$$\min_{\pi \in \Pi} \max_{\beta \in \Delta(\Sigma(\mathcal{B}))} L(\pi, \beta). \quad (8)$$

241 One can also prove that this above game has a value. Additionally, a solution (π_R^*, β_R^*) solving (8)
 242 has similar properties to (6) and (7) with respect to regret: π_R^* has its Bayesian regret upper bounded
 243 by $L(\pi_R^*, \beta_R^*)$ on $\Sigma(\mathcal{B})$ and is optimal on β_R^* .

244 Should utility or regret be used as an objective? Exploiting Regret ensures that scenarios from which
 245 you can learn the most from are sampled more often. It also ensures that degenerate scenarios get
 246 discarded as their regret is always zero. However, it demands the calculation of best responses for
 247 each scenario, which becomes taxing as the number of scenarios or problem complexity grows. To
 248 reduce the computational burden, we can approximate those best responses, or subsample the set of
 249 scenarios. An alternative way is to make use of the utility notion under some additional conditions:

250 **Definition 1** (Non-degenerative population). A background population of policies $\mathcal{B} \subset \Pi$ is non-
 251 degenerative $\iff \forall \sigma \in \Sigma(\mathcal{B}), \exists \pi_1, \pi_2 \in \Pi, \pi_1 \neq \pi_2$ and $U(\pi_1, \sigma) \neq U(\pi_2, \sigma)$.

252 **Lemma 1.** If a population \mathcal{B} is non-degenerative, then $\forall \sigma \in \Sigma(\mathcal{B}), \exists \pi \in \Pi, R(\pi, \sigma) > 0$.

253 *Proof.* \mathcal{B} is non-degenerative, for any scenario $\sigma \in \Sigma(\mathcal{B})$ there must exist two policies π_1 and π_2
 254 such that $U(\pi_1, \sigma) > U(\pi_2, \sigma)$. We have by definition $U^*(\sigma) \geq U(\pi_1, \sigma)$, hence $R(\pi_2, \sigma) > 0$. \square

255 Making the assumption that a background population is non-degenerative is in general realistic for
 256 cooperative tasks. This translates into only considering reasonable behaviors for the background
 257 population, or tasks where teammates cannot completely cancel out the actions of the focal players.
 258 Under the assumption of a non-degenerative population, no distribution can deadlock the policy
 259 learner into a stale scenario. For the remainder of the paper, background populations are assumed to
 260 be non-degenerative.

261 6 Computing Solutions

262 We desire to calculate the solution pairs for both the maximin utility (4) and minimax regret (8) games.
 263 Buening et al. (2023) theoretically proved that GDA has convergence guarantees when the game is
 264 played between a policy learned with softmax parameterization and nature learning its distribution
 265 over a finite set of MDPs. These results apply if all scenarios induce single-agent POMGs, as partial
 266 observability does not interfere with proving the required properties. However, when the focal policy
 267 is deployed in a scenario with $c > 1$ copies, the game is no longer single-agent.

268 To approximate the reduction of these multi-agent POMGs to single-agent POMGs during training,
 269 we propose using delayed versions π_{t-d} of the focal policy π_t for the $c - 1$ remaining copies. This
 270 common practice smooths the behavior of the copies and favors proper convergence by treating the
 271 copies as fixed policies. An implementation of GDA for our setting is provided in Appendix A.

272 7 Experiments

273 The aim of our experiments is to highlight the importance of partner distribution in the learning
 274 process. To achieve this, we evaluate our proposed strategies, Maximin Utility (MU) and Minimax
 275 Regret (MR), on two distinct problems. First, we consider the fully known and observable repeated
 276 Prisoner’s Dilemma to validate the theoretical results. Following this, we test our approaches on
 277 a deep-learning task, the Collaborative Cooking (Overcooked) game (Carroll et al., 2019; Strouse
 278 et al., 2021; Leibo et al., 2021; Agapiou et al., 2023). Throughout our experiments, we benchmark
 279 MU and MR against other distribution management strategies: SP which fixes the prior as the Dirac
 280 distribution $\beta^{\text{SP}}(\sigma^{\text{SP}}) = 1$, FP which is similar to SP but has the versions of the copies sampled
 281 uniformly from the full history of policies π_0, \dots, π_t , and Population Best-Response (PBR) which
 282 learns the best response to the training background population by maintaining a uniform prior
 283 $\beta^{\text{PBR}} = \mathcal{U}(\Sigma(\mathcal{B}^{\text{train}}))$. Approaches are consistently evaluated on their training background population,
 284 as well as on a separate test set, in order to evaluate their AHT capabilities.

285 7.1 Repeated Prisoner’s Dilemma

286 In these experiments, all computations can be exact. This includes the gradient calculation for the
 287 prior, as well as for the agent’s policies. We focus on the repeated Prisoner’s Dilemma, where two
 288 players play the matrix game repeatedly for $T = 3$ rounds.

Table 1: Scores on the repeated Prisoner’s Dilemma. A higher value is desired for performance (p , average utility) and worst-case utility (U^-), while a lower worst-case regret (R^+) is better. $p(\beta)$ corresponds to the utility w.r.t. a specific distribution β , rather than the average as in (3).

	$\Sigma(\mathcal{B}^{\text{train}})$					$\Sigma(\mathcal{B}^{\text{test}})$		
	p	$p(\beta_U^*)$	$p(\beta_R^*)$	U^-	R^+	p	U^-	R^+
MU	6.82	3.00	8.10	3.00	8.92	8.65	3.08	8.92
MR	9.14	0.92	11.25	0.92	2.11	7.54	0.99	8.45
PBR	10.0	1.96	10.66	1.96	3.00	7.74	1.99	10.46
FP	9.69	0.14	11.89	0.14	2.95	7.10	0.17	10.53
SP	9.69	0.46	11.99	0.46	3.00	7.21	0.47	10.70
TfT	10.0	2.00	11.73	2.00	3.00	7.60	2.01	10.65
CuD	10.0	2.00	11.73	2.00	3.00	7.85	2.01	9.92
Random	8.10	1.50	10.10	1.5	4.5	8.00	2.52	4.5

289 **Experimental Setup.** In the repeated Prisoner’s Dilemma, players receive and observe rewards based
 290 on their chosen actions: players receive a reward of 1 if both defect, 4 if both cooperate, 5 and 0 if
 291 the first defects while the second cooperates. The game has one state, and the outcomes observed are
 292 enough to determine the joint actions, making it fully observable.

293 We use softmax, fully adaptive policies, where actions depend on the entire history of observations
 294 and actions. During training, the learner has access to a background population containing a pure
 295 cooperative policy, a pure defective policy, and two popular strategies for the game: Tit-for-Tat (TfT),
 296 which mimics the partner’s previous action and starts with a cooperate action, and Cooperate-until-
 297 Defected (CuD), which defects if any defection was observed in the past, otherwise cooperating. A
 298 separate test population $\mathcal{B}^{\text{test}}$ is generated beforehand by randomly sampling 32 stochastic policies.

299 **Results.** Table 1 presents scores on the training and test sets. The results on the training set confirm
 300 most expectations: PBR is the best under the uniform prior (p), MU has the highest worst-case utility
 301 (U^-), and MR has the lowest worst-case regret (R^+). However, MR is not optimal on the worst-case
 302 prior β_R^* , likely due to the approximate self-play. On the test set, the best performance is achieved by
 303 MU, rather than PBR. MU also has the highest worst-case utility. Lastly, the random strategy excels
 304 in terms of worst-case regret on the test set, likely because it avoids fully committing to defecting or
 305 cooperating. Besides the random strategy, MR has the most respectable worst-case regret, closely
 306 followed by MU. These results indicate that learning the best response to populations does not ensure
 307 the best robustness to new partners. We also remark that SP and FP agents seem to overfit to their
 308 own established conventions, resulting in poor transferability to training and test policies.

309 7.2 Robust Cooperation in Deep RL Tasks

310 For this section, we tackle the Collaborative Cooking game (Agapiou et al., 2023), where two players
 311 act as chefs in a gridworld kitchen, working together to deliver as many tomato soup dishes as
 312 possible within a set time. To do so, they have to collect tomatoes, cook them, prepare dishes, and
 313 deliver the soup. Successful deliveries reward both players equally. Players must navigate the kitchen,
 314 interact with objects in the right order, and coordinate with each other. Each player has an egocentric,
 315 partial RGB view of the environment. All of our policies in this section are using deep recurrent
 316 (LSTM) neural networks.

317 **Experimental Setup.** Two separate background populations, $\mathcal{B}^{\text{train}}$ and $\mathcal{B}^{\text{test}}$, are generated according
 318 to Section 4.1. Both populations are trained with an identical setup, differing only in their seed. Each
 319 is partitioned into four sub-populations of sizes 2, 3, and 5, totaling 10 policies. Prosociality and
 320 risk-aversion are sampled uniformly in $[-0.2, 1.2]$ and $[0.1, 2]$ respectively. The same populations
 321 are used throughout three random seeds.

322 For a fair comparison and to focus on scenario distribution learning, we assume that $\mathcal{B}^{\text{train}}$ is readily
 323 available to all approaches, which can be exploited for a maximum of 4×10^7 environment steps to
 324 learn a policy with PPO (Schulman et al., 2017). We evaluate the approaches on two different kitchen
 325 layouts: Circuit and Cramped (Agapiou et al., 2023).

Table 2: Scores on the Collaborative Cooking environment training set. The standard error is taken over three random seeds. The scores are aggregated over kitchen layouts.

	$\Sigma(\mathcal{B}^{\text{train}})$				
	p	$p(\beta_U^*)$	$p(\beta_R^*)$	U^-	R^+
MU	266.9 ± 4.3	23.1 ± 0.7	24.5 ± 0.8	225.3 ± 11.5	266.0 ± 7.9
MR	232.0 ± 18.6	19.9 ± 1.6	20.2 ± 1.4	144.3 ± 14.4	230.7 ± 28.2
PBR	209.7 ± 23.9	20.0 ± 1.7	18.4 ± 3.1	96.8 ± 13.4	357.6 ± 16.1
FP	129.9 ± 13.9	7.9 ± 0.7	12.0 ± 1.1	0.2 ± 0.1	483.5 ± 16.1
SP	124.8 ± 26.4	9.4 ± 2.8	11.8 ± 3.2	15.7 ± 10.5	460.8 ± 21.7
Random	42.8 ± 0.0	6.7 ± 0.0	4.2 ± 0.0	0.0 ± 0.0	505.4 ± 0.0

Table 3: Scores on the Collaborative Cooking environment test sets.

	$\Sigma(\mathcal{B}^{\text{test}})$			$\Sigma^{\text{Melting Pot}}$		
	p	U^-	R^+	p	U^-	R^+
MU	195.7 ± 6.2	66.0 ± 6.8	266.4 ± 10.1	273.8 ± 4.9	224.9 ± 7.1	118.0 ± 7.1
MR	172.2 ± 15.4	65.1 ± 16.0	248.2 ± 28.4	206.8 ± 12.6	148.7 ± 9.1	187.1 ± 13.0
PBR	151.4 ± 19.5	33.6 ± 5.9	327.1 ± 14.0	171.8 ± 21.1	106.3 ± 9.3	228.1 ± 5.8
FP	152.2 ± 16.8	16.7 ± 11.7	369.7 ± 19.7	121.5 ± 15.7	40.7 ± 16.3	294.8 ± 10.7
SP	117.4 ± 12.5	6.7 ± 3.5	367.5 ± 11.6	101.2 ± 17.8	29.0 ± 17.4	293.4 ± 14.5
PP-ACB	n/a	n/a	n/a	82.4 ± 0.0	0.0 ± 0.0	307.3 ± 0.0
PP-OPRE	n/a	n/a	n/a	102.3 ± 0.0	14.6 ± 0.0	292.7 ± 0.0
PP-VMPO	n/a	n/a	n/a	78.6 ± 0.0	36.1 ± 0.0	306.7 ± 0.0
Random	32.2 ± 0.0	0.0 ± 0.0	445.0 ± 0.0	60.6 ± 0.0	0.0 ± 0.0	307.3 ± 0.0

326 Additionally, we assess the learned policies on the Melting Pot benchmark scenarios, comparing them
 327 against the scores of the baselines reported in the original paper (Agapiou et al., 2023): an Actor-Critic
 328 Baseline (ACB), V-MPO (Song et al., 2019), and OPRE (Vezhnevets et al., 2020). Note that these
 329 three baselines were trained through PP without our background policies, for 10^9 environment steps.

330 **Results.** The results in Table 2 and Table 3 clearly show that MU marginally outperforms all other
 331 evaluated methods. Looking at the robustness metrics, MU has the best worst-case utilities (U^-) and
 332 the best worst-case regret (R^+) on the Melting Pot scenarios. MR also performs better than any other
 333 benchmarked method overall, securing the lowest worst-case regrets on both the training and test sets.
 334 In terms of performance (p), MU and MR are consistently the best and second-best, respectively,
 335 which is particularly notable on the training set where PBR was expected to perform the best.

336 One hypothesis for MR performing worse than MU globally is that the estimations of the maximal
 337 utilities for training scenarios are too approximate. Another hypothesis for why MU and MR
 338 marginally outperform PBR and other approaches is that the distributions learned during training
 339 have a similar effect to curriculum learning, introducing indirect exploration in behaviours compared
 340 to fixed distributions.

341 8 Conclusion

342 We investigated how to obtain robust adaptive policies in an AHT setting. Leveraging work on
 343 Minimax-Bayes RL, we proposed a method to find worst-case distributions over background popu-
 344 lations, which led to consistently robust policies compared to simply training policies on uniform
 345 distributions. In addition, we have found the unexpected results that these distributional choices
 346 significantly accelerate learning. For future work, we believe that adapting our methods for a partner
 347 distribution-based curriculum-learning could be highly promising. This approach has the potential to
 348 strongly enhance sample efficiency, asymptotic performance, and robustness.

349 **References**

- 350 John P. Agapiou, Alexander Sasha Vezhnevets, Edgar A. Duéñez-Guzmán, Jayd Matyas, Yiran Mao,
351 Peter Sunehag, Raphael Köster, Udari Madhushani, Kavya Kopparapu, Ramona Comanescu, DJ
352 Strouse, Michael B. Johanson, Sukhdeep Singh, Julia Haas, Igor Mordatch, Dean Mobbs, and
353 Joel Z. Leibo. 2023. Melting Pot 2.0. arXiv:2211.13746 [cs.MA]
- 354 Stefano Albrecht, Jacob Crandall, and Subramanian Ramamoorthy. 2015. An Empirical Study on
355 the Practical Impact of Prior Beliefs over Policy Types. *Proceedings of the AAAI Conference on*
356 *Artificial Intelligence* 29, 1 (Feb. 2015).
- 357 Samuel Barrett, Avi Rosenfeld, Sarit Kraus, and Peter Stone. 2017. Making friends on the fly:
358 Cooperating with new teammates. *Artificial Intelligence* 242 (2017), 132–171.
- 359 Samuel Barrett and Peter Stone. 2015. Cooperating with Unknown Teammates in Complex Domains:
360 A Robot Soccer Case Study of Ad Hoc Teamwork. *Proceedings of the AAAI Conference on*
361 *Artificial Intelligence* 29, 1 (Feb. 2015).
- 362 Samuel Barrett, Peter Stone, and Sarit Kraus. 2011. Empirical evaluation of ad hoc teamwork in
363 the pursuit domain. In *The 10th International Conference on Autonomous Agents and Multiagent*
364 *Systems-Volume 2*. 567–574.
- 365 James O Berger. 1985. Statistical decision theory and Bayesian analysis. *Springer Series in Statistics*
366 (1985).
- 367 George W. Brown. 1951. Iterative Solution of Games by Fictitious Play. In *Activity Analysis of*
368 *Production and Allocation*, T. C. Koopmans (Ed.). Wiley, New York.
- 369 Thomas Kleine Buening, Christos Dimitrakakis, Hannes Eriksson, Divya Grover, and Emilio Jorge.
370 2023. Minimax-Bayes Reinforcement Learning. In *Proceedings of The 26th International Confer-*
371 *ence on Artificial Intelligence and Statistics (Proceedings of Machine Learning Research, Vol. 206)*,
372 Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent (Eds.). PMLR, 7511–7527.
- 373 Rodrigo Canaan, Xianbo Gao, Julian Togelius, Andy Nealen, and Stefan Menzel. 2023. Generating
374 and Adapting to Diverse Ad Hoc Partners in Hanabi. *IEEE Transactions on Games* 15, 2 (2023),
375 228–241.
- 376 Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca
377 Dragan. 2019. On the utility of learning about humans for human-ai coordination. *Advances in*
378 *neural information processing systems* 32 (2019).
- 379 Rujikorn Charakorn, Poramate Manoonpong, and Nat Dilokthanakul. 2020. Investigating partner
380 diversification methods in cooperative multi-agent deep reinforcement learning. In *Neural Infor-*
381 *mation Processing: 27th International Conference, ICONIP 2020, Bangkok, Thailand, November*
382 *18–22, 2020, Proceedings, Part V* 27. Springer, 395–402.
- 383 Shuo Chen, Ewa Andrejczuk, Zhiguang Cao, and Jie Zhang. 2020. Aateam: Achieving the ad hoc
384 teamwork by employing the attention mechanism. In *Proceedings of the AAAI conference on*
385 *artificial intelligence*, Vol. 34. 7095–7102.
- 386 Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. 2019. Quantifying
387 Generalization in Reinforcement Learning. In *Proceedings of the 36th International Conference*
388 *on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri
389 and Ruslan Salakhutdinov (Eds.). PMLR, 1282–1289.
- 390 Michael Dennis, Natasha Jaques, Eugene Vinitzky, Alexandre Bayen, Stuart Russell, Andrew Critch,
391 and Sergey Levine. 2020. Emergent complexity and zero-shot transfer via unsupervised environ-
392 ment design. *Advances in neural information processing systems* 33 (2020), 13049–13061.
- 393 Samuel Garcin, James Doran, Shangmin Guo, Christopher G Lucas, and Stefano V Albrecht. 2023.
394 How the level sampling process impacts zero-shot generalisation in deep reinforcement learning.
395 *arXiv preprint arXiv:2310.03494* (2023).

- 396 Adam Gleave, Michael Dennis, Cody Wild, Neel Kant, Sergey Levine, and Stuart Russell. 2019.
397 Adversarial policies: Attacking deep reinforcement learning. *arXiv preprint arXiv:1905.10615*
398 (2019).
- 399 Johannes Heinrich, Marc Lanctot, and David Silver. 2015. Fictitious Self-Play in Extensive-Form
400 Games. In *Proceedings of the 32nd International Conference on Machine Learning (Proceedings*
401 *of Machine Learning Research, Vol. 37)*, Francis Bach and David Blei (Eds.). PMLR, Lille, France,
402 805–813.
- 403 Irina Higgins, Arka Pal, Andrei Rusu, Loic Matthey, Christopher Burgess, Alexander Pritzel, Matthew
404 Botvinick, Charles Blundell, and Alexander Lerchner. 2017. DARLA: Improving Zero-Shot
405 Transfer in Reinforcement Learning. In *Proceedings of the 34th International Conference on*
406 *Machine Learning (Proceedings of Machine Learning Research, Vol. 70)*, Doina Precup and
407 Yee Whye Teh (Eds.). PMLR, 1480–1490.
- 408 Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob Foerster. 2020. “Other-Play” for Zero-
409 Shot Coordination. In *Proceedings of the 37th International Conference on Machine Learning*
410 *(Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.).
411 PMLR, 4399–4410.
- 412 Max Jaderberg, Wojciech M. Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia Cas-
413 tañeda, Charles Beattie, Neil C. Rabinowitz, Ari S. Morcos, Avraham Ruderman, Nicolas Sonnerat,
414 Tim Green, Louise Deason, Joel Z. Leibo, David Silver, Demis Hassabis, Koray Kavukcuoglu, and
415 Thore Graepel. 2019. Human-level performance in 3D multiplayer games with population-based
416 reinforcement learning. *Science* 364, 6443 (2019), 859–865.
- 417 Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan,
418 Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. 2023. Ai alignment: A comprehensive survey.
419 *arXiv preprint arXiv:2310.19852* (2023).
- 420 Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. 2021. Prioritized Level Replay. In *Proceed-*
421 *ings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning*
422 *Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 4940–4950.
- 423 Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Perolat,
424 David Silver, and Thore Graepel. 2017. A Unified Game-Theoretic Approach to Multiagent
425 Reinforcement Learning. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von
426 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30.
427 Curran Associates, Inc.
- 428 Joel Z Leibo, Edgar A Dueñez-Guzman, Alexander Vezhnevets, John P Agapiou, Peter Sunehag,
429 Raphael Koster, Jayd Matyas, Charlie Beattie, Igor Mordatch, and Thore Graepel. 2021. Scalable
430 Evaluation of Multi-Agent Reinforcement Learning with Melting Pot. In *Proceedings of the 38th*
431 *International Conference on Machine Learning (Proceedings of Machine Learning Research,*
432 *Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 6187–6199.
- 433 Shihui Li, Yi Wu, Xinyue Cui, Honghua Dong, Fei Fang, and Stuart Russell. 2019. Robust Multi-
434 Agent Reinforcement Learning via Minimax Deep Deterministic Policy Gradient. *Proceedings of*
435 *the AAAI Conference on Artificial Intelligence* 33, 01 (Jul. 2019), 4213–4220.
- 436 Tianyi Lin, Chi Jin, and Michael Jordan. 2020. On Gradient Descent Ascent for Nonconvex-Concave
437 Minimax Problems. In *Proceedings of the 37th International Conference on Machine Learning*
438 *(Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.).
439 PMLR, 6083–6093.
- 440 Andrei Lupu, Brandon Cui, Hengyuan Hu, and Jakob Foerster. 2021. Trajectory Diversity for Zero-
441 Shot Coordination. In *Proceedings of the 38th International Conference on Machine Learning*
442 *(Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.).
443 PMLR, 7204–7213.
- 444 Reuth Mirsky, Ignacio Carlucho, Arrasy Rahman, Elliot Fosong, William Macke, Mohan Sridharan,
445 Peter Stone, and Stefano V. Albrecht. 2022. A Survey of Ad Hoc Teamwork Research. In *Multi-*
446 *Agent Systems*, Dorothea Baumeister and Jörg Rothe (Eds.). Springer International Publishing,
447 Cham, 275–293.

- 448 Darius Muglich, Christian Schroeder de Witt, Elise van der Pol, Shimon Whiteson, and Jakob Foerster.
449 2022a. Equivariant Networks for Zero-Shot Coordination. In *Advances in Neural Information*
450 *Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.),
451 Vol. 35. Curran Associates, Inc., 6410–6423.
- 452 Darius Muglich, Luisa M Zintgraf, Christian A Schroeder De Witt, Shimon Whiteson, and Jakob
453 Foerster. 2022b. Generalized Beliefs for Cooperative AI. In *Proceedings of the 39th International*
454 *Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika
455 Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (Eds.).
456 PMLR, 16062–16082.
- 457 Frans A Oliehoek. 2012. Decentralized pomdps. In *Reinforcement learning: state-of-the-art*. Springer,
458 471–503.
- 459 Alexander Peysakhovich and Adam Lerer. 2017. Prosocial learning agents solve generalized stag
460 hunts better than selfish ones. *arXiv preprint arXiv:1709.02865* (2017).
- 461 Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. 2017. Robust Adversarial
462 Reinforcement Learning. In *Proceedings of the 34th International Conference on Machine Learning*
463 *(Proceedings of Machine Learning Research, Vol. 70)*, Doina Precup and Yee Whye Teh (Eds.).
464 PMLR, 2817–2826.
- 465 Manish Ravula, Shani Alkoby, and Peter Stone. 2019. Ad Hoc Teamwork With Behavior Switch-
466 ing Agents. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial*
467 *Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization,
468 550–556.
- 469 Michael Rovatsos and Marco Wolf. 2002. Towards social complexity reduction in multiagent learning:
470 the adhoc approach. In *Proceedings of the 2002 AAAI Spring Symposium on Collaborative Learning*
471 *Agents*. 90–97.
- 472 Lukas Schäfer. 2022. Task generalisation in multi-agent reinforcement learning. In *Proceedings of*
473 *the 21st International Conference on Autonomous Agents and Multiagent Systems*. 1863–1865.
- 474 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal
475 policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- 476 H Francis Song, Abbas Abdolmaleki, Jost Tobias Springenberg, Aidan Clark, Hubert Soyer, Jack W
477 Rae, Seb Noury, Arun Ahuja, Siqi Liu, Dhruva Tirumala, et al. 2019. V-mpo: On-policy maximum a
478 posteriori policy optimization for discrete and continuous control. *arXiv preprint arXiv:1909.12238*
479 (2019).
- 480 Peter Stone, Gal Kaminka, Sarit Kraus, and Jeffrey Rosenschein. 2010. Ad hoc autonomous agent
481 teams: Collaboration without pre-coordination. In *Proceedings of the AAAI Conference on Artificial*
482 *Intelligence*, Vol. 24. 1504–1509.
- 483 DJ Strouse, Kevin McKee, Matt Botvinick, Edward Hughes, and Richard Everett. 2021. Collaborating
484 with humans without human data. *Advances in Neural Information Processing Systems* 34 (2021),
485 14502–14515.
- 486 Alexander Vezhnevets, Yuhuai Wu, Maria Eckstein, Rémi Leblond, and Joel Z Leibo. 2020. Options
487 as REsponses: Grounding behavioural hierarchies in multi-agent reinforcement learning. In
488 *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine*
489 *Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, 9733–9742.

Algorithm 1 Background-Focal GDA

- 1: **Input** set of background policies \mathcal{B} , and learning rates (η_π, η_β) .
 - 2: Initialize randomly the main policy parameters θ_0
 - 3: Initialize the belief as the uniform distribution over possible scenarios $\beta_0 = \mathcal{U}(\Sigma(\mathcal{B}))$
 - 4: **for** $t = 0, \dots, N - 1$ **do**
 - 5: Compute $U(\pi_{\theta_t}, \sigma)$ for all $\sigma \in \Sigma(\mathcal{B})$
 - 6: Compute $U(\pi_{\theta_t}, \beta_t) = \sum_{\sigma} U(\pi_{\theta_t}, \sigma) \beta_t(\sigma)$
 - 7: Compute $R(\pi_{\theta_t}, \sigma) = U^*(\sigma) - U(\pi_{\theta_t}, \sigma)$ for all $\sigma \in \Sigma(\mathcal{B})$
 - 8: Obtain $L(\pi_{\theta_t}, \beta_t) = \sum_{\sigma} R(\pi_{\theta_t}, \sigma) \beta_t(\sigma)$
 - 9: Update belief $\beta_{t+1} = \mathcal{P}(\beta_t + \eta_\beta \nabla_{\beta} L(\pi_{\theta_t}, \beta_t))$ (projection onto the simplex)
 - 10: Update policy parameters $\theta_{t+1} = \theta_t + \eta_\theta \nabla_{\theta} U(\pi_{\theta_t}, \beta_t)$
 - 11: **end for**
 - 12: **return** θ^*, β^* uniformly at random from $\{(\theta_1, \beta_1), \dots, (\theta_N, \beta_N)\}$
-

Algorithm 2 Background-Focal SGDA

- 1: **Input** set of background policies \mathcal{B} , batch size B , learning rates (η_π, η_β)
 - 2: Initialize randomly the main policy parameters θ_0
 - 3: Initialize the belief as the uniform distribution over possible scenarios $\beta_0 = \mathcal{U}(\Sigma(\mathcal{B}))$
 - 4: **for** $t = 0, \dots, N - 1$ **do**
 - 5: Sample B scenarios $\sigma_1, \dots, \sigma_B \sim \beta_t$
 - 6: Estimate $\hat{U}(\pi_{\theta_t}, \sigma_i)$ by deploying π_{θ_t} on σ_i , for $i = 1, \dots, B$
 - 7: Compute $\hat{U}(\pi_{\theta_t}, \beta_t) = \frac{1}{B} \sum_{i=1}^B \hat{U}(\pi_{\theta_t}, \sigma_i)$
 - 8: Compute $\hat{R}(\pi_{\theta_t}, \sigma_i) = U^*(\sigma_i) - \hat{U}(\pi_{\theta_t}, \sigma_i)$ for each $i = 1, \dots, B$
 - 9: Compute $\hat{L}(\pi_{\theta_t}, \beta_t) = \frac{1}{B} \sum_{i=1}^B \hat{R}(\pi_{\theta_t}, \sigma_i)$
 - 10: Update belief $\beta_{t+1} = \mathcal{P}(\beta_t + \eta_\beta \nabla_{\beta} \hat{L}(\pi_{\theta_t}, \beta_t))$ (projection onto the simplex)
 - 11: Update policy parameters $\theta_{t+1} = \theta_t + \eta_\theta \nabla_{\theta} \hat{U}(\pi_{\theta_t}, \beta_t)$
 - 12: **end for**
 - 13: **return** θ^*, β^* uniformly at random from $\{(\theta_1, \beta_1), \dots, (\theta_N, \beta_N)\}$
-

490 **Appendix**491 **A Algorithms**

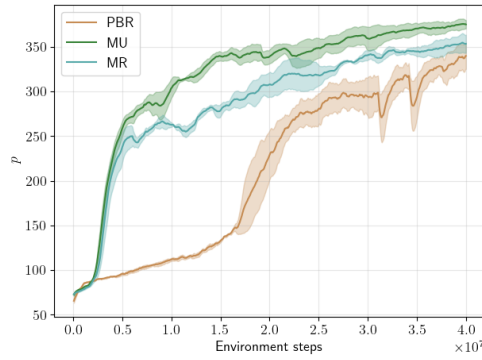
492 We provide implementations for solving the minimax regret problem with (Algorithm 1) or without
 493 (Algorithm 2) full knowledge of the game. Updating the belief with the rule

$$\beta_{t+1} = \mathcal{P}(\beta_t - \eta_\beta \nabla_{\beta} U(\pi_{\theta_t}, \beta_t)) \quad (9)$$

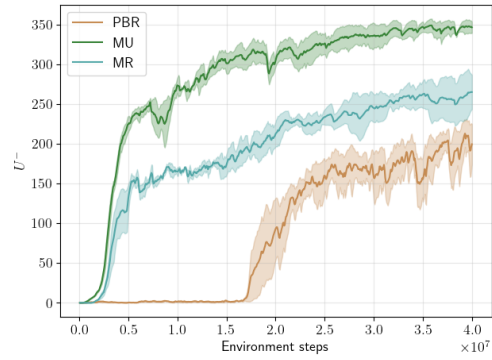
494 solves the maximin utility problem instead.

495 **B Additional experimental results**496 **B.1 Robust Cooperation in Deep RL Tasks**

497 We provide learning curves for the Collaborative Cooking game. For both kitchen layouts, the curves
 498 in Figures 2 and 3 uncover the fact that both the minimax regret and maximin utility formulations
 499 significantly speed up learning. The prior curves provided in Figures 4 and 5, highly smoothed for
 500 interpretability, show how different distributions are learned with utility and regret.

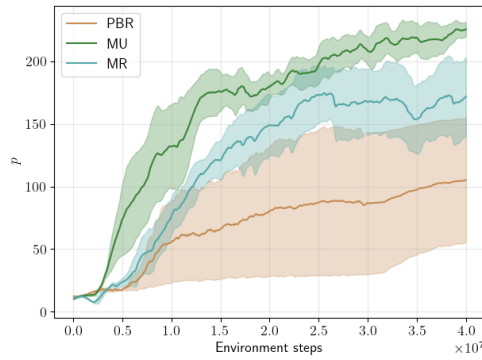


(a) Performance.

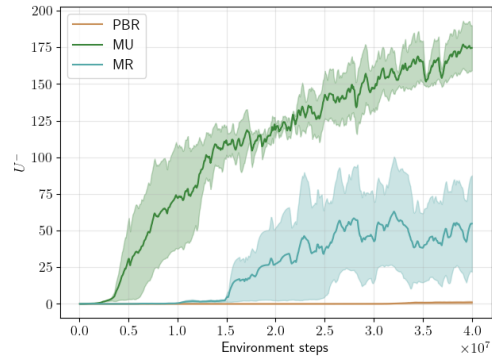


(b) Worst-case utility.

Figure 2: Learning curves of the average and worst-case utility metrics over the training set of the Collaborative Cooking Cramped environment. The standard error is taken over three random seeds.

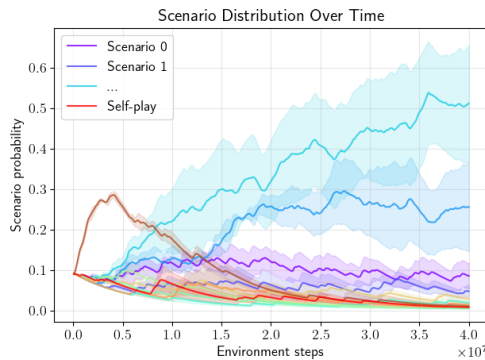


(a) Performance.

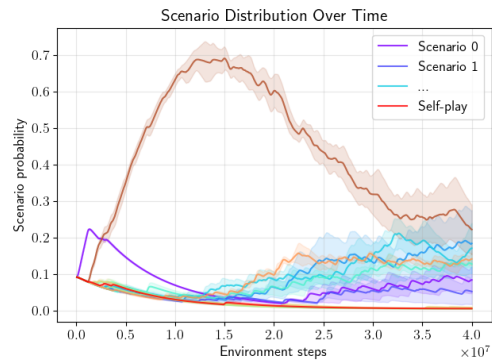


(b) Worst-case utility.

Figure 3: Learning curves of the average and worst-case utility metrics over the training set of the Collaborative Cooking Circuit environment. The standard error is taken over three random seeds.



(a) Maximin Utility.



(b) Minimax Regret.

Figure 4: Learning curves of the prior over the training scenarios, for the Collaborative Cooking Cramped environment. The standard error is taken over three random seeds.

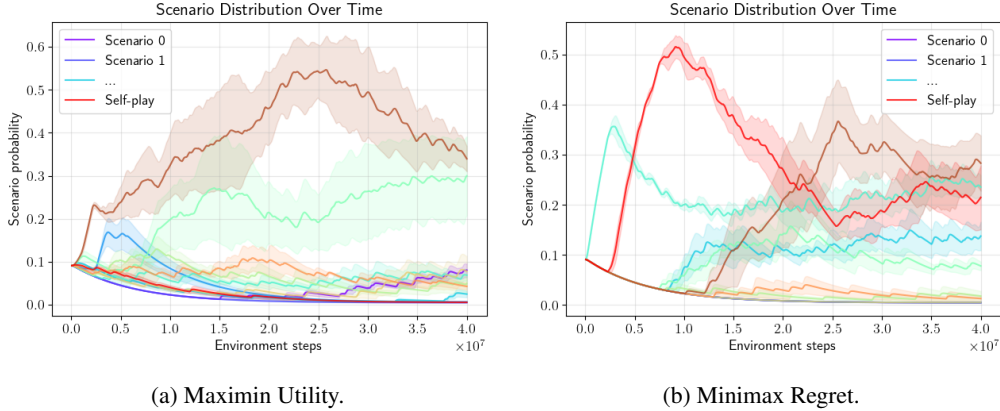


Figure 5: Learning curves of the prior over the training scenarios, for the Collaborative Cooking Circuit environment. The standard error is taken over three random seeds.

501 C Additional experimental details

502 C.1 Robust Cooperation in Deep RL Tasks

503 To facilitate the training of our policies in Collaborative Cooking, we used a shaping pseudo-reward
 504 of 1 when tomatoes were placed in the cooking pot. For the background policies, we further altered
 505 the reward function to restrict delivery rewards to the player that delivered. Combining this new
 506 reward function with varying levels of prosociality and risk-aversion helped the background policies
 507 adopt diversified ways to solve the game.

508 The architecture for the agents consisted of a convolutional network with two layers, having 16 and
 509 32 output channels, kernel shapes of 8 and 4, and strides of 8 and 1, respectively. The output of the
 510 convNet was concatenated with the previous action taken before being passed into a dense layer of
 511 size 256 and an LSTM with 256 units. Policy and baseline (for the critic) were produced by linear
 512 layers connected to the output of the LSTM.

513 We chose PPO to train our policies, using the Adam optimizer with a learning rate of 2×10^{-4} ,
 514 a discount factor of 0.99, a GAE lambda of 0.95, a KL coefficient of 1.0 with a KL target of
 515 0.01, and a PPO clipping parameter of 0.3. Gradients were clipped at 4.0. We did not employ
 516 entropy regularization. PPO was set to run 2 epochs per batch, each containing 64000 samples, with
 517 minibatches of 1000 samples each. Finally, the unroll length for the LSTM was set at 20.

518 For the prior, we used a learning rate of 0.4. We also constrained the prior to keep a probability of
 519 $5e - 2$ to sample a random scenario in order to allow constant exploration.