

---

# Molten Pot: Evaluations & Datasets for Social Offline Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many of the settings where reinforcement learning (RL) could matter most are both  
2 social and data-limited: agents must act in the presence of other decision-makers,  
3 yet cannot rely on online interaction to learn how to do so. Current evaluations do  
4 not target data-constrained social reasoning. Standard offline RL benchmarks treat  
5 the environment as non-social, while multi-agent benchmarks focus exclusively on  
6 fully cooperative settings. Thus the challenge of reasoning about partner identity  
7 and motivations under mixed incentives, and generalising across social structures  
8 from offline data alone, remains untested. We introduce Molten Pot, an evaluation  
9 protocol, datasets and benchmark for offline mixed-motive social RL built on  
10 Melting Pot substrates. The protocol spans five substrates, 47 social scenarios,  
11 approximately one terabyte of trajectory data, and defines three complementary  
12 evaluation settings that each probe a different aspect of social robustness. Setting  
13 1 tests offline RL in multi-agent, mixed-motive settings with fixed background  
14 populations. Setting 2 pools datasets across every scenario of a substrate, requiring  
15 the learnt policy to handle varying partner behaviour without explicit context.  
16 Setting 3 evaluates zero-shot social generalisation through disjoint train/test splits  
17 that isolate specific social shifts. Finally, we benchmark four representative offline  
18 RL algorithms on our evaluation protocol and datasets, finding clear limitations  
19 in current methods’ ability to learn robust social strategies from offline data alone.  
20 Molten Pot establishes offline social evaluation as a distinct and necessary target  
21 for RL research.

## 22 1 Introduction

23 Humans routinely act appropriately in mixed-motive social situations they have never encountered  
24 before [4, 39]. A driver merging into unfamiliar traffic, a negotiator facing an unknown counterpart,  
25 a statistician designing an adaptive trial for a patient population whose adherence patterns must be  
26 estimated: each must infer the intentions and likely behaviour of others from limited observation  
27 and adapt without the luxury of *online* trial and error [3]. This capacity for social reasoning in  
28 novel settings is crucial for trusting agents to act in the real world [9]. Yet in many of the situations  
29 where such agents would be deployed, abundant interaction data already exists while further online  
30 learning is expensive, dangerous, or impractical. Consider autonomous vehicles trained on recorded  
31 driving data alongside human drivers whose behaviour varies and cannot be controlled. Alternatively,  
32 consider designing adaptive or platform trials from a corpus of past trials run by different sponsors  
33 and sites, whose enrollment, dosing, and early-stopping practices vary across the dataset and cannot  
34 be intervened on after the fact [17]. In each case the data is plentiful, the environment is social, and  
35 the agent has no opportunity to explore online [17, 27]. We must evaluate social reasoning under the  
36 data constraints that deployment frequently imposes. Current evaluation infrastructure does not test  
37 this capability.

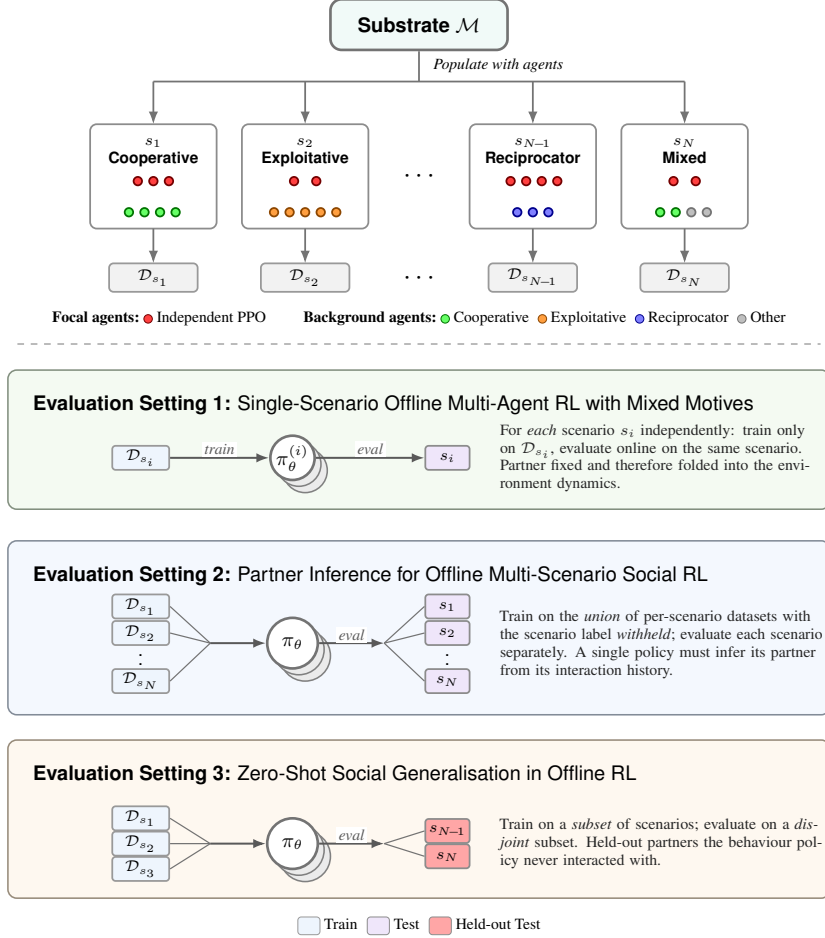


Figure 1: **Molten Pot at a glance.** *Top:* a substrate  $\mathcal{M}$  paired with different background populations  $\pi_s^B$  defines a family of scenarios  $s_i = \langle \mathcal{M}, \pi_{s_i}^B \rangle$ . Every scenario hosts multiple *focal* agents, each controlled by its own independent PPO learner; trajectories were logged throughout training so the offline dataset  $\mathcal{D}_{s_i}$  mixes skill levels from random initialisation to the converged policy. *Bottom:* the same datasets feed three offline RL evaluation settings that probe distinct facets of social robustness. Evaluation Setting 1 trains and evaluates per scenario independently; Evaluation Setting 2 pools all datasets with the scenario label withheld and tests partner inference; Evaluation Setting 3 trains on a subset of scenarios and evaluates on a disjoint held-out set of scenarios.

38 Offline RL, as benchmarked by D4RL [14], RL Unplugged [19], and their successors, have focused  
 39 almost exclusively on single-agent control problems, in which the environment dynamics are sta-  
 40 tionary and the training distribution varies only through the quality of the behaviour policy used  
 41 to collect data. The evaluation targets these benchmarks stress are distributional shift between the  
 42 behaviour policy  $\pi^\beta$  and the learned policy  $\pi^\theta$ , and the conservatism required to avoid over-estimating  
 43 out-of-distribution actions. What they do *not* evaluate are the social reasoning challenges that arise  
 44 when the environment contains other adaptive agents. In a social environment, the standard offline  
 45 RL question of whether we can learn a robust policy from historical logs, becomes whether we can  
 46 learn a robust policy from logs of interactions between a policy and a *population of co-players with*  
 47 *mixed motives*, and does that *policy generalize to unseen players*?

48 Offline multi-agent RL benchmarks move closer but do not reach this goal. OG-MARL and its  
 49 extensions [10, 12] target offline Multi-Agent RL (MARL) in fully cooperative environments such as  
 50 SMAC [36] and the MAMuJoCo [33], where the central challenge is coordination [38] rather than  
 51 handling heterogeneous partner behaviour under misaligned incentives. Mixed-motive settings, in  
 52 which agents have overlapping but non-identical incentives and must reason about whether and when  
 53 to cooperate, are largely absent. Melting Pot [25] provides exactly the substrates needed to probe

54 mixed-motive social reasoning, but was designed for online evaluation: policies are trained with  
 55 self-play or population-based methods and then scored against held-out background populations. The  
 56 offline counterpart — learning mixed-motive social policies from fixed datasets and evaluating on  
 57 these substrates — remains unexplored.

58 We introduce *Molten Pot*, an evaluation protocol with datasets for offline mixed-motive, social  
 59 reinforcement learning. It is built on five Melting Pot substrates (Figure 2) that together span several  
 60 social reasoning challenges including: public-goods provision, agent reciprocity, partner choice,  
 61 commons governance, and trust formation. Molten Pot provides infrastructure to evaluate offline  
 62 mixed-motive social reinforcement learning: per-scenario offline datasets totalling approximately one  
 63 terabyte of trajectory data, and three complementary evaluation settings that probe distinct aspects of  
 64 social robustness (Figure 1).

65 *Evaluation Setting 1* tests single-scenario offline learning with fixed partner populations, mirroring  
 66 the standard offline MARL setup [10] but now testing the ability to navigate a mixed-motive social  
 67 environment. *Evaluation Setting 2* pools datasets across every scenario of a substrate and withholds the  
 68 scenario label, requiring policies to handle heterogeneous partner behaviour without explicit context,  
 69 whether by explicit partner inference, by learning a robust generalist strategy, or otherwise. *Evaluation*  
 70 *Setting 3* designs train and test splits with held-out scenarios that exhibit specific classes of social  
 71 shift, including: convention mismatch, reciprocity threshold, punishment capability, deployment role,  
 72 and latent within-episode partner type. This isolates each type of shift and tests whether methods can  
 73 generalize to partners exhibiting behaviours the training data never contained. The three settings are  
 74 not arranged on a difficulty ladder; rather they probe different challenges and a method may succeed  
 75 on one and fail on another for reasons specific to that setting.

76 Finally, the capability Molten Pot targets is increasingly relevant beyond the classical RL community.  
 77 Large language model agents are being deployed into open-ended environments [22] where they  
 78 interact with other agents (other LLM agents, humans, and legacy automated systems) whose policies  
 79 are unknown and must be inferred from the example interactions seen during pre-training. Molten Pot  
 80 provides precisely the infrastructure to ask "*how well do these agents handle unfamiliar co-players*  
 81 *given only its fixed pre-training data?*". We view the benchmark as a precursor to the systematic  
 82 evaluation of deployed LLM-based agents under similar constraints.

## 83 2 Preliminaries

### 84 2.1 Partially Observable Markov Games

85 Following Leibo et al. [25], a *substrate* is a multi-agent environment defined by its map, mechanics,  
 86 and reward structure, independent of the specific players that populate it. Formally, a substrate is  
 87 modelled as a partially observable general-sum Markov game [18, 20]

$$\mathcal{M} = \langle \mathcal{N}, \mathcal{Z}, \{\mathcal{A}_i\}_{i \in \mathcal{N}}, \{\mathcal{O}_i\}_{i \in \mathcal{N}}, \mathcal{T}, \{r_i\}_{i \in \mathcal{N}}, \gamma \rangle,$$

88 with players  $\mathcal{N}$ , state space  $\mathcal{Z}$ , per-player discrete action and observation spaces  $\mathcal{A}_i, \mathcal{O}_i$ , joint  
 89 transition kernel  $\mathcal{T} : \mathcal{Z} \times \prod_i \mathcal{A}_i \rightarrow \Delta(\mathcal{Z})$ , per-player reward  $r_i : \mathcal{Z} \times \prod_i \mathcal{A}_i \rightarrow \mathbb{R}$ , and discount  
 90  $\gamma \in [0, 1)$ . Here, observations are egocentric  $88 \times 88$  RGB frames and rewards are individual;  
 91 games are mixed-motive, meaning reward functions are neither identical (cooperative) nor opposed  
 92 (zero-sum) unless explicitly stated.

### 93 2.2 Focal and Background Agents

94 Following Leibo et al. [25], the player set is partitioned into *focal* positions  $\mathcal{F} \subset \mathcal{N}$ , whose policies  
 95 are learned and evaluated, and *background* positions  $\mathcal{B} = \mathcal{N} \setminus \mathcal{F}$ , whose policies are fixed and shipped  
 96 with the benchmark. Within the benchmarked implementations, focal agents share parameters  $\theta$   
 97 across positions; background agents are scripted or pre-trained co-players that act as a black box  
 98 from the learner’s perspective. Parameter sharing is one choice for tackling focal learning.

### 99 2.3 Scenarios

100 A *scenario* is a pair  $s = \langle \mathcal{M}_s, \pi_s^{\mathcal{B}} \rangle$  comprising a substrate  $\mathcal{M}_s$  and a fixed joint background policy  
 101  $\pi_s^{\mathcal{B}}$ . Molten Pot ships a finite set  $\mathcal{S}_{\text{all}}$  of scenarios across five substrates (Tables 2–3); each scenario

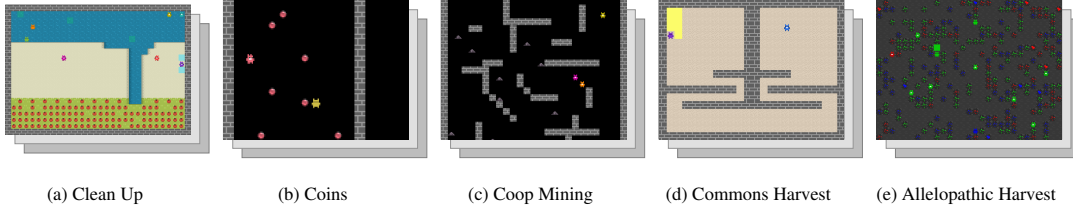


Figure 2: **The five Molten Pot substrates.** Representative rendered frames from DeepMind’s Melting Pot [25] covering a range of mixed-motive game structures: common-pool resource dilemmas, public-goods games, dyadic social dilemmas, partner-choice coordination, and convention adoption. Every substrate ships multiple background-partner scenarios: Clean Up (9), Coins (7), Coop Mining (6), Commons Harvest (12), and Allelopathic Harvest (13), for a total of 47 scenarios. Per-scenario descriptions appear in Table 3.

102 captures a qualitatively distinct social configuration, such as a focal population *visiting* a resident  
 103 group of reciprocators or *being visited* by a corrigible rider (one that cooperates only after being  
 104 sanctioned). Fixing a scenario  $s$  absorbs the background agents into the environment dynamics:  
 105 from the focal perspective, background actions  $\mathbf{a}^B \sim \pi_s^B(\cdot | h_t^B)$  are part of the transition function  
 106 rather than strategic variables. Writing  $\mathcal{M}_s^F(\pi_\theta)$  for the resulting focal-perspective dynamics under  
 107 scenario  $s$  and focal policy  $\pi_\theta$ , the focal return is

$$J_s(\pi_\theta) = \mathbb{E}_{\tau \sim \mathcal{M}_s^F(\pi_\theta)} \left[ \sum_{t=0}^{T-1} \gamma^t r_t^F \right], \quad r_t^F = \frac{1}{|\mathcal{F}|} \sum_{i \in \mathcal{F}} r_{i,t}. \quad (1)$$

## 108 2.4 Datasets

109 For a scenario  $s$ , an offline dataset  $\mathcal{D}_s$  is a collection of focal-perspective episodes generated by  
 110 some focal agent policies acting in  $\mathcal{M}_s$  alongside the fixed background population  $\pi_s^B$ . Each episode  
 111 records, for every focal agent  $i \in \mathcal{F}$  and timestep  $t$ , the observation  $o_{i,t}$ , action  $a_{i,t}$ , and reward  
 112  $r_{i,t}$ . Background actions and the underlying environment state are *not* stored: from the learner’s  
 113 perspective, partner behaviour is observable only through its effect on the focal trajectory.

114 For every scenario  $s$ , the focal agents are trained against the fixed background population  $\pi_s^B$  using  
 115 *independent* PPO [37]: each focal agent runs its own PPO learner and updates on its own individual  
 116 reward, so from any single focal agent’s perspective the learning problem is a partially observable  
 117 stochastic game [18] in which the background is fixed and the other focal agents are concurrent,  
 118 non-stationary co-learners. This inherits the canonical difficulties of multi-agent RL; non-stationarity  
 119 from simultaneous focal updates and credit assignment over shared environmental events.

120 During each run we log  $N = 1000$  trajectories of length  $T = 1000$  spaced uniformly across training  
 121 time rather than only at convergence, so  $\mathcal{D}_s$  deliberately mixes focal-skill levels, from near-random  
 122 early-training behaviour through to the converged policy. To validate the quality of our converged poli-  
 123 cies, the focal agent returns are compared to the returns reported by Leibo et al. [25] on the same sce-  
 124 narios, confirming that the behaviour policy is a reasonable ceiling against which to normalise offline  
 125 methods. Following the dataset-transparency recommendation of Formanek et al. [12], the aggregated  
 126 return distribution per substrate is shown in Figure 3 and the per-scenario breakdown in Appendix B.

## 127 3 Partner Inference and Bayes Optimality.

128 Unlike Evaluation Setting 1, Settings 2 and 3 both require focal agents to act well against varying  
 129 partner populations that are never identified by label. Setting 2 is designed specifically so that  
 130 maximising its evaluation metric formally requires partner inference. We now show why.

131 Because background behaviour is unobserved and varies across scenarios, the dynamics experienced  
 132 by a focal agent depend on a *latent context*: the background partners themselves. This latent-context  
 133 structure mirrors the social inference problem faced by humans entering unfamiliar social settings:  
 134 observed behaviour reveals information about the intentions and dispositions of others, which in turn  
 135 informs action selection [3]. We can formalise this structure. Writing  $h_t = (o_0, a_0, r_0, \dots, o_t)$  for

136 the focal history, the focal-perspective transition under scenario  $s$  is obtained by marginalising over  
 137 background actions:

$$p(o_{t+1}, r_t | h_t, a_t; s) = \mathbb{E}_{\mathbf{a}^B \sim \pi_s^B(\cdot | h_t^B)} \left[ p(o_{t+1}, r_t | h_t, a_t, \mathbf{a}^B) \right]. \quad (2)$$

138 When focal agents are deployed against a scenario drawn from a set  $\mathcal{S}$  with prior  $\Pr(s)$ , the Bayes-  
 139 filter posterior over scenarios given the observed history is

$$b_t(s) := \Pr(s | h_t) \propto \Pr(s) \prod_{u=0}^{t-1} p(o_{u+1}, r_u | h_u, a_u; s). \quad (3)$$

140 As focal agents accumulate observations and rewards, each timestep provides evidence about which  
 141 partner population it faces, and the belief  $b_t$  concentrates on the true scenario. The Bayes-optimal focal  
 142 policy acts under the belief-weighted predictive dynamics  $p(o_{t+1}, r_t | h_t, a_t) = \sum_s b_t(s) p(o_{t+1}, r_t |$   
 143  $h_t, a_t; s)$ , and the value of the optimal such policy satisfies the Bellman equation

$$V^*(h_t) = \max_{a_t} \mathbb{E}_{s \sim b_t} \left[ \mathbb{E}_{r_t, o_{t+1} \sim p(\cdot | h_t, a_t; s)} [r_t + \gamma V^*(h_{t+1})] \right]. \quad (4)$$

144 This is the optimal expected return for the scenario distribution  $\mathcal{S}$  under prior  $\Pr(s)$ : because  $V^*$   
 145 solves the maximisation over all history-conditioned policies given exact knowledge of each scenario’s  
 146 generative model and exact Bayesian inference, no other history-conditioned policy achieves higher  
 147 expected return averaged over scenario draws from the same prior. On any individual realised scenario  
 148 a policy better calibrated to that specific partner population may achieve higher return, but only by  
 149 sacrificing performance on others;  $V^*$  is the optimum over the distribution as a whole, and serves as  
 150 an upper bound on the expected performance of any learned policy, where both the generative model  
 151 and the posterior must additionally be approximated from finite offline data. The tension between  
 152 per-scenario specialisation and distributional robustness is precisely what the evaluation protocol is  
 153 designed to expose.

154 This Bayesian formulation characterises what must be solved, not how. An agent that matches  
 155 the Bayes-optimal policy through explicit posterior tracking and an agent that arrives at the same  
 156 value through a learned generalist strategy are equally optimal; in the latter case, the requisite social  
 157 reasoning is occurring implicitly. Maximising expected return averaged over a scenario distribution,  
 158 where partner identity is unobserved and inferable only from within-episode history, therefore  
 159 formally requires partner inference, whether performed explicitly or implicitly. The evaluation  
 160 protocol below, in particular Setting 2, operationalises this directly.

## 161 4 The Molten Pot Evaluation Protocol

162 Molten Pot defines three *evaluation settings* that probe different aspects of social robustness in offline  
 163 RL. Each is a fully specified offline reinforcement-learning problem; what distinguishes them is the  
 164 support of the training and evaluation scenario distributions. The five substrates, the full scenario  
 165 catalogue, and the train/test splits of evaluation setting 3 are documented in Appendix A (Tables 2–4).

166 Formally, we write  $\mathcal{D}_{\mathcal{S}} := \bigcup_{s \in \mathcal{S}} \mathcal{D}_s$  for the union of the per-scenario datasets over a set of scenarios  
 167  $\mathcal{S} \subseteq \mathcal{S}_{\text{all}}$ . Each evaluation setting specifies a *training* scenario set  $\mathcal{S}^{\text{tr}}$ , whose datasets the learner  
 168 may use, and a *test* scenario set  $\mathcal{S}^{\text{te}}$ , in which the learned focal policy  $\pi_{\theta}$  is scored. The settings  
 169 differ in how these two sets relate: identical in Settings 1 and 2, disjoint in Setting 3. A learner is  
 170 given  $\mathcal{D}_{\mathcal{S}^{\text{tr}}}$  and must return  $\pi_{\theta}$  without any environment interaction. In every setting the scenario  
 171 label  $s$  is withheld both from the observation at training time and from the policy at test time: the  
 172 algorithm only ever sees focal  $(o, a, r)$  tuples. Our datasets are released in full on our anonymised  
 173 Hugging Face repository<sup>1</sup>.

### 174 4.1 Evaluation Setting 1: Single-Scenario Offline Multi-Agent RL with Mixed Motives

175 **Setup.** Setting 1 is run independently per scenario: for a given  $s \in \mathcal{S}_{\text{all}}$ , we set  $\mathcal{S}^{\text{tr}} = \mathcal{S}^{\text{te}} = \{s\}$ .  
 176 The algorithm is trained on  $\mathcal{D}_s$  and evaluated online in  $\mathcal{M}_s$ . Performance is the focal return from (1)  
 177 evaluated in scenario  $s$ :

$$J_{\text{I}}(\pi_{\theta}; s) = J_s(\pi_{\theta}), \quad s \in \mathcal{S}_{\text{all}}. \quad (5)$$

<sup>1</sup><https://huggingface.co/datasets/frmjua/moltenpot> (Anonymous)

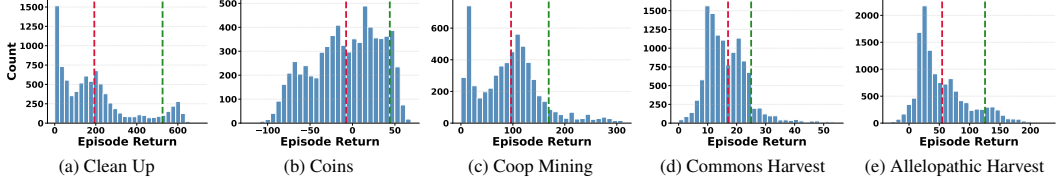


Figure 3: **Per-substrate focal-agent return distributions.** Each panel pools all trajectories logged during the independent-PPO runs of every scenario in the substrate. The spread of episode returns confirms that every substrate’s dataset is skill-mixed rather than expert-only. The dashed red line in each panel marks the substrate’s mean episode return, and the dashed green line marks the 90th percentile. Across the benchmark we logged  $\sim 1$  TB of focal-agent trajectories. Per-scenario return breakdowns appear in Appendix B.

178 **What this setting evaluates.** Because  $s$  is fixed,  $\pi_s^{\mathcal{B}}$  is a fixed feature of the environment dynamics  
 179 and no partner inference is required. Setting 1 therefore mirrors the standard single-task offline POSG  
 180 from Formanek et al. [10], but in a mixed-motive multi-agent environment instead of a cooperative  
 181 one. This highlights a gap in existing evaluation: offline RL benchmarks do not test mixed-motive  
 182 social environments, and multi-agent offline benchmarks are restricted to fully cooperative settings  
 183 [10, 11]. Setting 1 fills that gap directly, asking whether standard offline methods can learn to act  
 184 well alongside partners who do not share rewards.

#### 185 4.2 Evaluation Setting 2: Partner Inference for Offline Multi-Scenario Social RL

186 **Setup.** Setting 2 is run independently per substrate: for a given substrate, let  $\mathcal{S}_{\text{sub}} \subseteq \mathcal{S}_{\text{all}}$  denote  
 187 its scenarios. We set  $\mathcal{S}^{\text{tr}} = \mathcal{S}^{\text{te}} = \mathcal{S}_{\text{sub}}$  and train a *single* focal policy on the union dataset  $\mathcal{D}_{\mathcal{S}_{\text{sub}}}$ .  
 188 Crucially, the scenario label is *not* provided to the learner. The trained policy is then deployed and  
 189 evaluated in each scenario separately. The headline score is the mixture return

$$J_{\text{II}}(\pi_{\theta}) = \mathbb{E}_{s \sim P^{\text{te}}(\mathcal{S}_{\text{sub}})} [J_s(\pi_{\theta})], \quad (6)$$

190 where  $P^{\text{te}}$  is the uniform distribution over  $\mathcal{S}_{\text{sub}}$ . We also report the per-scenario profile  
 191  $\{J_s(\pi_{\theta})\}_{s \in \mathcal{S}_{\text{sub}}}$  so that scenario specific failures are visible (Fig. 5).

192 **What this setting evaluates.** Setting 2 instantiates the evaluation in §3: because evaluation episodes  
 193 are independent, the policy is deployed fresh in each scenario and must infer the partner population  
 194 from within-episode history alone. The headline metric  $J_{\text{II}}$  averages across scenarios under the  
 195 uniform prior, so its maximiser is the Bayes-optimal partner-inferring policy. A policy that ignores  
 196 partner identity and plays a single marginal strategy is sub-optimal. The inverse cannot be claimed:  
 197 poor performance need not stem from missing partner inference alone. High performance also  
 198 requires generative model specification from finite offline data, and standard offline-RL errors such  
 199 as extrapolation error, insufficient support, and conservative bias can also drive poor performance.  
 200 No existing benchmark provides a setting in which partner inference capability can be measured.

#### 201 4.3 Evaluation Setting 3: Zero-Shot Social Generalisation in Offline RL

202 **Setup.** Setting 3 is run independently per substrate. For each substrate we provide disjoint train/test  
 203 splits  $(\mathcal{S}^{\text{tr}}, \mathcal{S}^{\text{te}})$  with  $\mathcal{S}^{\text{tr}} \cap \mathcal{S}^{\text{te}} = \emptyset$  (Table 4). The algorithm is trained on  $\mathcal{D}_{\mathcal{S}^{\text{tr}}}$ , again with scenario  
 204 labels withheld, and evaluated on the held-out scenarios. Performance is  $J_{\text{II}}$  as in Setting 2, but  
 205 evaluated on the held-out scenarios  $\mathcal{S}^{\text{te}}$ .

206 **Formal view.** Setting 3 introduces a distribution shift  $P^{\text{tr}}(s) \neq P^{\text{te}}(s)$  with disjoint scenario  
 207 support: the partner posterior (Eq. 3) is *miscalibrated by construction*, and the Bayes-optimal rule (4)  
 208 cannot be evaluated directly. The setting is nonetheless not designed to be unsolvable: the observation  
 209 space, action space, and substrate mechanics are shared across all three settings, so the distributional  
 210 shift is in the composition of familiar elements rather than in the elements themselves.

211 **What this setting evaluates.** Setting 3 introduces the evaluation of zero-shot social generalisation  
 212 from offline data. Where Setting 2 evaluates partner inference over a known scenario set, Setting 3

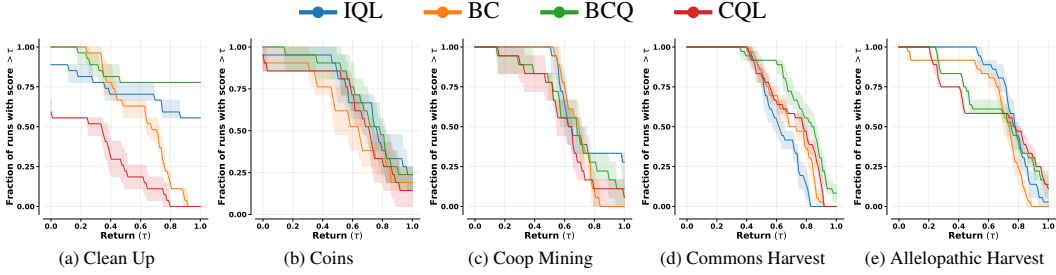


Figure 4: **Evaluation 1 — per-substrate performance profiles.** For each substrate we plot the reliable performance profile [2, 16]: the fraction of scenarios on which each algorithm scores at least the normalised final-eval return given on the  $x$ -axis. Shaded bands are stratified-bootstrap CIs.

213 asks whether learned social reasoning transfers to partner populations entirely absent from the training  
 214 distribution. Each train/test split isolates a specific, nameable social shift: convention mismatch,  
 215 threshold shift in reciprocity, punishment-capability shift, role shift, commitment shift, latent within-  
 216 episode partner type, feedback-loop shift, and population-scale shift. Table 4 in Appendix A lists the  
 217 nineteen splits across the five substrates and characterises each shift in detail.

## 218 5 Benchmark Results

219 We report results on Evaluation 1, Evaluation 2, and Evaluation 3 using four representative algorithms:  
 220 behavioural cloning [BC, 35], Batch-Constrained Q-Learning [BCQ, 15], Implicit Q-Learning [IQL,  
 221 23], and Conservative Q-Learning [CQL, 24]. Every algorithm’s policy and Q-function carries a  
 222 recurrent GRU [6] layer over the focal agent’s interaction history so that partner identity can in  
 223 principle be inferred from past observations and actions [30]. We use single-agent offline algorithms  
 224 because existing offline MARL methods target fully cooperative settings via centralised training,  
 225 which is ill-suited to mixed-motive substrates with distinct per-agent rewards; independent-agent  
 226 baselines have moreover been shown to be competitive with, and often to outperform, CTDE methods  
 227 on standardised offline MARL benchmarks [13]. Each algorithm is trained for 20,000 gradient  
 228 updates with three seeds, and hyperparameters are included in Appendix C. Our code is available on  
 229 GitHub<sup>2</sup>.

230 **Normalisation.** For each scenario we report the normalised final eval return  
 231  $\tilde{J}_s = (R_s - R_{\text{random},s}) / (R_{p90,s} - R_{\text{random},s})$  [14], where  $R_{\text{random},s}$  is the mean return of a  
 232 uniform-random focal policy on scenario  $s$  and  $R_{p90,s}$  is the 90th-percentile episode return in the  
 233 offline dataset for  $s$  (which we use as a proxy for expert policies). Under this normalisation 0  
 234 corresponds to a random policy and 1 to the “best-in-dataset” (top-decile) episode. The reference  
 235 lines on every result figure mark  $\tilde{J} = 0$  (grey dotted, random) and  $\tilde{J} = 1$  (green dashed, dataset  $p_{90}$ ).

### 236 5.1 Evaluation Setting 1 — Aggregate Performance across Scenarios

237 Figure 4 reports a *performance profile* per substrate: the fraction of the substrate’s scenarios on  
 238 which each algorithm achieves at least the normalised final-eval return on the  $x$ -axis. Curves are  
 239 computed by MARL-evals utilities [16], with shaded stratified-bootstrap confidence intervals [2]. On  
 240 most substrates BC produces the most consistent profile — a near-flat plateau followed by a sharp  
 241 drop at a substrate-specific threshold, indicating that BC scores in a narrow band around roughly  
 242 the same return on most scenarios. The offline RL methods (IQL, BCQ, CQL) trace more gradually  
 243 decreasing curves, reflecting wider per-scenario spread: their curves sit *below* BC’s at thresholds  
 244 inside BC’s plateau (scenarios where they underperform BC) but retain non-zero mass to the right of  
 245 BC’s drop-off (scenarios where they substantially exceed BC). This reflects the offline RL algorithms’  
 246 sensitivity to the varying magnitudes of expected episode return across scenarios: sometimes they  
 247 fail to learn anything useful, and sometimes they learn a policy that outperforms the dataset’s  $p_{90}$   
 248 episodes by a large margin. Per-scenario breakdowns appear in Appendix D. In Appendix F we also  
 249 include the change in background agents’ returns after the focal agent’s policy is deployed. In a

<sup>2</sup><https://github.com/frmjua/moltenpot/> (Anonymous)

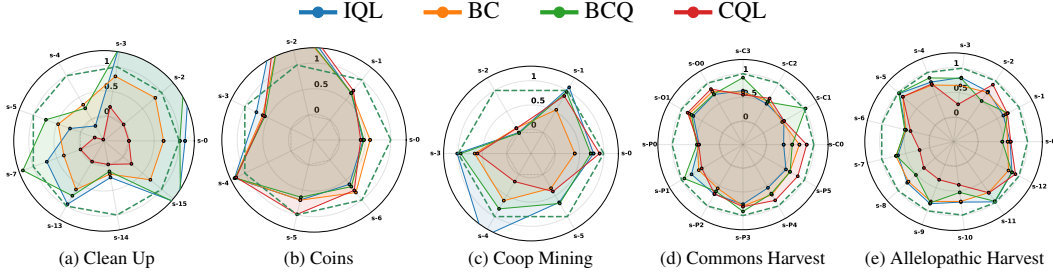


Figure 5: **Evaluation 2 — per-substrate scenario profile.** The rays of each spider are that substrate’s scenarios (numbered along the perimeter, with Commons Harvest variants abbreviated C=closed, O=open, P=partnership); each algorithm’s polygon connects its mean normalised final return on each scenario. The dashed green circle marks the dataset  $p_{90}$  ceiling and the dotted grey circle marks the random-policy baseline. Asymmetric polygons reveal scenarios on which an algorithm collapses.

250 few scenarios, especially Coop Mining, the background agents returns increase along with the focal  
 251 agents, reflecting a *win-win* situation. But in the majority of scenarios, the focal agent’s gain is at the  
 252 expense of the background agents, reflecting the mixed-motive nature of the substrates.

## 253 5.2 Evaluation Setting 2 — Per-Substrate Profiles

254 The substrate-level aggregate hides whether an algorithm is uniformly mediocre or instead does  
 255 well on a few scenarios and collapses on others. Figure 5 unpacks Setting 2 by substrate, with one  
 256 spider plot per substrate whose rays are scenarios and whose polygon vertices are each algorithm’s  
 257 mean normalised return. The polygons reveal within-substrate variation that the aggregate erases: on  
 258 every substrate at least one scenario lies at or below the random-policy reference for every algorithm,  
 259 and the dent toward the centre falls in broadly the same place across algorithms—consistent with  
 260 the posterior-inference reading, in which the difficult scenarios are those whose partner type is  
 261 mismatched against the training mixture rather than the algorithm’s choice of conservatism penalty.  
 262 Per-scenario Setting 1 vs Setting 2 comparisons appear in Appendix D.

263 **Setting 1 vs Setting 2.** Pooling every (substrate, scenario, seed) sample, Table 1 reports RLi-  
 264 able’s [2, 16] optimality-gap headline estimator with stratified-bootstrap 95% confidence intervals. Despite Setting 2 training on  
 265 the union of all per-scenario datasets — substantially more data, drawn from a more diverse mixture than any single  
 266 Setting 1 run — the optimality gap tends to grow. We tentatively attribute this pattern to the partner-inference  
 267 pressure that Setting 2 imposes on the learner: withholding the scenario label forces the policy to recover the partner  
 268 posterior from interaction history alone, and the offline algorithms studied here appear unable to convert the extra  
 269 data into better generalist social behaviour.  
 270  
 271  
 272  
 273  
 274  
 275

Table 1: **Optimality gap** ( $\downarrow$ ) pooled across every (substrate, scenario, seed), with rliable 95% CI. All aggregate metrics in Appendix G.

Algo	Setting 1	Setting 2
IQL	0.30 [0.28, 0.33]	0.31 [0.29, 0.32]
BC	0.35 [0.33, 0.36]	0.36 [0.35, 0.37]
BCQ	<b>0.26 [0.24, 0.27]</b>	<b>0.29 [0.27, 0.31]</b>
CQL	0.43 [0.42, 0.44]	0.49 [0.48, 0.51]

## 276 5.3 Evaluation Setting 3 — Train vs. Test Performance.

277 Figure 6 reports each algorithm’s normalised return aggregated across (split, scenario, seed) samples  
 278 within each substrate, with in-distribution (train, filled circle) and out-of-distribution (held-out test,  
 279 filled square) evaluations shown side by side. Within every substrate the in-distribution markers sit  
 280 above their out-of-distribution counterparts: the disjoint-scenario shift costs a measurable fraction of  
 281 return for every algorithm, and on several substrates pushes out-of-distribution performance back  
 282 towards the random-policy reference. The gap size varies by algorithm and substrate. Punishment-  
 283 and role-shifts on Commons Harvest produce the largest drops, consistent with the hypothesis  
 284 that capability shifts (passive  $\rightarrow$  sanctioning counterparts; visitor  $\rightarrow$  resident deployment) require  
 285 composing two transfer problems at once. Coins absorbs its splits most gracefully, where train  
 286 and test scenarios share a common dyadic structure and differ mainly in partner thresholds. No

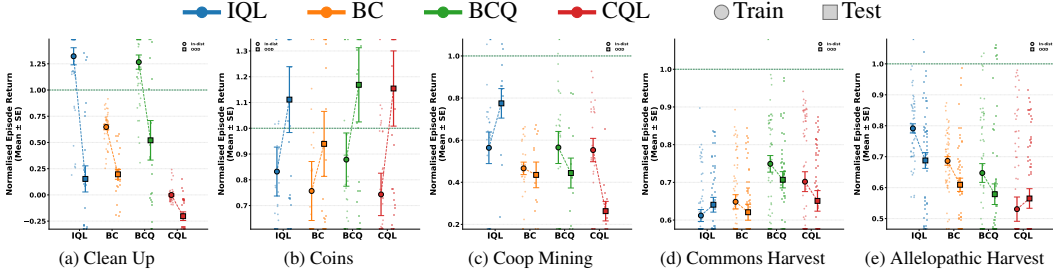


Figure 6: **Evaluation 3 — per-substrate aggregate normalised return, in-distribution vs out-of-distribution.** For each algorithm: markers are means pooled across every (split, scenario, seed) sample; error bars are SE. Outlier samples beyond the clipped axes are shown as pointing triangles at the edge. Per-split, per-scenario breakdowns in Appendix E.

287 algorithm closes the in/out gap: scenario-agnostic offline RL recovers within-distribution behaviour  
 288 but does not produce policies that generalise compositionally to disjoint partner populations. Per-split,  
 289 per-scenario raw-return breakdowns appear in Appendix E.

## 290 6 Related Work

291 Standard single-agent offline RL benchmarks—D4RL [14] and RL Unplugged [19]—fix per-task  
 292 behaviour distributions so that conservative algorithms such as BCQ [15], CQL [24], and IQL [23]  
 293 can be compared on a common axis of distributional shift, but their non-social dynamics make the  
 294 failure modes Molten Pot targets structurally invisible. The closest multi-agent counterparts are  
 295 OG-MARL [10] and its revision [12], which inherit the cooperative assumption from SMAC [36] and  
 296 MAMuJoCo [33] and therefore do not evaluate partner reasoning; offline-MARL algorithmic work  
 297 has likewise concentrated on cooperative coordination [5, 21, 38], multi-agent conservatism [32],  
 298 model-based approaches [5, 41], and offline-to-online transfer [28]. The mixed-motive evaluation  
 299 protocol Molten Pot adopts splits generalisation to unseen co-players, and which has since been  
 300 reimplemented for fast online training in SocialJax [29]; in the mixed-motive online setting, Ndousse  
 301 et al. [31] further show that social observation can drive agents to learn from one another. Partner  
 302 handling itself is evaluated in cooperative form by ZSC-Eval [40] and the ad-hoc-teamwork literature,  
 303 with Mon-Williams et al. [30] showing that recurrent agents form partner models only when the task  
 304 demands them. Applied mixed-motive MARL has also begun to make its way into human-compatible  
 305 self-driving [7, 8]. Molten Pot inherits the offline-data discipline of OG-MARL [10, 12] while  
 306 replacing its cooperative assumption with the focal/background mixed-motive structure of Melting  
 307 Pot, isolating partner handling rather than coordination as the object of evaluation.  
 308

## 309 7 Conclusion

310 We introduced Molten Pot<sup>3</sup>, an offline RL evaluation protocol with datasets built from five mixed-  
 311 motive Melting Pot substrates. Our three Evaluation Settings probe distinct aspects of social robust-  
 312 ness: single-scenario fitting (Setting 1), partner inference under withheld scenario labels (Setting 2),  
 313 and zero-shot generalisation across disjoint train/test partner populations (Setting 3). Four representa-  
 314 tive algorithms (CQL, BC, BCQ, IQL) leave substantial headroom against the  $p_{90}$  dataset ceiling on  
 315 Settings 1 and 2 and incur a clear in-distribution to out-of-distribution drop on Setting 3. Together, the  
 316 three settings evaluate failure modes that standard offline RL benchmarks overlook. Tackling these  
 317 settings will require methods that go beyond scenario-agnostic conservatism—partner representations  
 318 that support extrapolation, posterior-aware conservatism penalties, or memory architectures designed  
 319 to maintain a calibrated belief over partner type. Molten Pot provides the datasets, splits, and evalua-  
 320 tion harness needed to measure progress on each axis, and we hope it serves as a foundation for the  
 321 next generation of multi-agent offline RL.

<sup>3</sup>Interactive visualization: <https://frmjua.github.io/moltenpot> (Anonymous)

## References

- [1] John P. Agapiou, Alexander Sasha Vezhnevets, Edgar A. Duéñez-Guzmán, Jayd Matyas, Yiran Mao, Peter Sunehag, Raphael Košter, Udari Madhushani, Kavya Kopparapu, Ramona Comanescu, et al. Melting pot 2.0, 2022.
- [2] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G. Belle-mare. Deep reinforcement learning at the edge of the statistical precipice. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. URL <https://arxiv.org/abs/2108.13264>.
- [3] Chris L. Baker, Rebecca Saxe, and Joshua B. Tenenbaum. Action understanding as inverse planning. *Cognition*, 113(3):329–349, 2009. ISSN 0010-0277. DOI: <https://doi.org/10.1016/j.cognition.2009.07.005>. URL <https://www.sciencedirect.com/science/article/pii/S0010027709001607>. Reinforcement learning and higher cognition.
- [4] Chris L. Baker, Julian Jara-Ettinger, Rebecca Saxe, and Joshua B. Tenenbaum. Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour*, 1, 2017. URL <https://api.semanticscholar.org/CorpusID:3338320>.
- [5] Paul Barde, Jakob Nicolaus Foerster, Derek Nowrouzezahrai, and Amy Zhang. A model-based solution to the offline multi-agent reinforcement learning coordination problem. *ArXiv*, abs/2305.17198, 2023. URL <https://api.semanticscholar.org/CorpusID:258959276>.
- [6] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014.
- [7] Daphne Cornelisse and Eugene Vinitsky. Human-compatible driving partners through data-regularized self-play reinforcement learning. *ArXiv*, abs/2403.19648, 2024. URL <https://api.semanticscholar.org/CorpusID:268732678>.
- [8] Marco Francis Cusumano-Towner, David Hafner, Alexander Hertzberg, Brody Huval, Aleksei Petrenko, Eugene Vinitsky, Erik Wijmans, Taylor W. Killian, Stuart Bowers, Ozan Sener, Philipp Kraehenbuehl, and Vladlen Koltun. Robust autonomy emerges from self-play. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=y0XoJpG6Qy>.
- [9] Allan Dafoe, Edward Hughes, Yoram Bachrach, Tantum Collins, Kevin R. McKee, Joel Z. Leibo, K. Larson, and Thore Graepel. Open problems in cooperative ai. *ArXiv*, abs/2012.08630, 2020. URL <https://api.semanticscholar.org/CorpusID:229220772>.
- [10] Claude Formanek, Asad Jeewa, Jonathan Shock, and Arnu Pretorius. Off-the-grid marl: Datasets and baselines for offline multi-agent reinforcement learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems, AAMAS '23*, pp. 2442–2444, Richland, SC, 2023. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9781450394321.
- [11] Juan Claude Formanek, Callum Rhys Tilbury, Louise Beyers, Jonathan Phillip Shock, and Arnu Pretorius. Dispelling the mirage of progress in offline MARL through standardised baselines and evaluation. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=CaAJeNkceP>.
- [12] Juan Claude Formanek, Louise Beyers, Callum Rhys Tilbury, Jonathan Phillip Shock, and Arnu Pretorius. Putting data at the centre of offline multi-agent reinforcement learning. *Journal of Data-centric Machine Learning Research*, 2026. URL <https://openreview.net/forum?id=Rp6H7FKkpf>.
- [13] Juan Claude Formanek et al. Dispelling the mirage of progress in offline MARL through standardised baselines and evaluation. 2025. Stub — please verify metadata.

- 371 [14] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: Datasets for  
372 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 373 [15] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning  
374 without exploration. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of*  
375 *the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine*  
376 *Learning Research*, pp. 2052–2062. PMLR, 09–15 Jun 2019. URL [https://proceedings.](https://proceedings.mlr.press/v97/fujimoto19a.html)  
377 [mlr.press/v97/fujimoto19a.html](https://proceedings.mlr.press/v97/fujimoto19a.html).
- 378 [16] Rihab Gorsane, Omayma Mahjoub, Ruan John de Kock, Roland Dubb, Siddarth Singh, and  
379 Arnu Pretorius. Towards a standardised performance evaluation protocol for cooperative marl.  
380 *Advances in Neural Information Processing Systems*, 35:5510–5521, 2022.
- 381 [17] Omer Gottesman, Fredrik D. Johansson, Joshua Meier, Jack Dent, Donghun Lee, Srivatsan  
382 Srinivasan, Linying Zhang, Yi Ding, David Wihl, Xuefeng Peng, Jiayu Yao, Isaac Lage,  
383 Christopher Mosch, Li wei H. Lehman, Matthieu Komorowski, A. Aldo Faisal, Leo An-  
384 thony Celi, David A. Sontag, and Finale Doshi-Velez. Evaluating reinforcement learning  
385 algorithms in observational health settings. *ArXiv*, abs/1805.12298, 2018. URL [https:](https://api.semanticscholar.org/CorpusID:44152464)  
386 [//api.semanticscholar.org/CorpusID:44152464](https://api.semanticscholar.org/CorpusID:44152464).
- 387 [18] Sven Gronauer and Klaus Diepold. Multi-agent deep reinforcement learning: a survey. *Artif. In-*  
388 *tell. Rev.*, 55(2):895–943, February 2022. ISSN 0269-2821. DOI: 10.1007/s10462-021-09996-w.  
389 URL <https://doi.org/10.1007/s10462-021-09996-w>.
- 390 [19] Caglar Gulcehre, Ziyu Wang, Alexander Novikov, Thomas Paine, Sergio Gómez, Konrad Zolna,  
391 Rishabh Agarwal, Josh Merel, Daniel Mankowitz, Cosmin Paduraru, Gabriel Dulac-Arnold,  
392 Jerry Li, Mohammad Norouzi, Matt Hoffman, Nicolas Heess, and Nando de Freitas. RL  
393 unplugged: A suite of benchmarks for offline reinforcement learning. In *Advances in Neural*  
394 *Information Processing Systems*, 2020.
- 395 [20] Eric A. Hansen, Daniel S. Bernstein, and Shlomo Zilberstein. Dynamic programming for  
396 partially observable stochastic games. In *Proceedings of the 19th National Conference on*  
397 *Artificial Intelligence*, AAAI’04, pp. 709–715. AAAI Press, 2004. ISBN 0262511835.
- 398 [21] Marcel Hedman, Kale ab Abebe Tessera, Juan Claude Formanek, Anya Sims, Riccardo Zamboni,  
399 Trevor McInroe, John Torr, and Elliot Fosong. Coda: Coordination via on-policy diffusion  
400 for multi-agent offline reinforcement learning, 2026. URL [https://arxiv.org/abs/2604.](https://arxiv.org/abs/2604.23308)  
401 [23308](https://arxiv.org/abs/2604.23308).
- 402 [22] Edward Hughes, Michael D Dennis, Jack Parker-Holder, Feryal Behbahani, Aditi Mavalankar,  
403 Yuge Shi, Tom Schaul, and Tim Rocktäschel. Position: Open-endedness is essential for artificial  
404 superhuman intelligence. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian  
405 Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the*  
406 *41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine*  
407 *Learning Research*, pp. 20597–20616. PMLR, 21–27 Jul 2024. URL [https://proceedings.](https://proceedings.mlr.press/v235/hughes24a.html)  
408 [mlr.press/v235/hughes24a.html](https://proceedings.mlr.press/v235/hughes24a.html).
- 409 [23] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit  
410 q-learning. *ArXiv*, abs/2110.06169, 2021. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:238634325)  
411 [CorpusID:238634325](https://api.semanticscholar.org/CorpusID:238634325).
- 412 [24] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning  
413 for offline reinforcement learning. *Advances in neural information processing systems*, 33:  
414 1179–1191, 2020.
- 415 [25] Joel Z. Leibo, Edgar Dué nez Guzmán, Alexander Sasha Vezhnevets, John P. Agapiou, Peter  
416 Sunehag, Raphael Koster, Jayd Matyas, Charles Beattie, Igor Mordatch, and Thore Graepel.  
417 Scalable evaluation of multi-agent reinforcement learning with melting pot. PMLR, 2021. DOI:  
418 10.48550/arXiv.2107.06857. URL <https://doi.org/10.48550/arXiv.2107.06857>.
- 419 [26] Adam Lerer and Alexander Peysakhovich. Maintaining cooperation in complex social dilemmas  
420 using deep reinforcement learning, 2018. URL <https://arxiv.org/abs/1707.01068>.

- 421 [27] Sergey Levine, Aviral Kumar, G. Tucker, and Justin Fu. Offline reinforcement learning:  
422 Tutorial, review, and perspectives on open problems. *ArXiv*, abs/2005.01643, 2020. URL  
423 <https://api.semanticscholar.org/CorpusID:218486979>.
- 424 [28] Trevor A. McInroe, Stefano V. Albrecht, and Amos J. Storkey. Planning to go out-of-distribution  
425 in offline-to-online reinforcement learning. *ArXiv*, abs/2310.05723, 2023. URL <https://api.semanticscholar.org/CorpusID:263829022>.  
426
- 427 [29] Ruoyu Mizuta et al. Socialjax: An evaluation suite for multi-agent reinforcement learning in  
428 sequential social dilemmas, 2025.
- 429 [30] Ruaridh Mon-Williams, Max Taylor-Davies, Elizabeth Mieczkowski, Natalia Velez, Neil R.  
430 Bramley, Yanwei Wang, Thomas L. Griffiths, and Christopher G. Lucas. Partner modelling  
431 emerges in recurrent agents (but only when it matters). *arXiv preprint arXiv:2505.17323*, 2025.  
432 URL <https://arxiv.org/abs/2505.17323>.
- 433 [31] Kamal Ndousse, Douglas Eck, Sergey Levine, and Natasha Jaques. Emergent social learning  
434 via multi-agent reinforcement learning. In *International Conference on Machine Learning*,  
435 2020. URL <https://api.semanticscholar.org/CorpusID:235621490>.
- 436 [32] Ling Pan, Longbo Huang, Tengyu Ma, and Huazhe Xu. Plan better amid conservatism: Offline  
437 multi-agent reinforcement learning with actor rectification. *CoRR*, abs/2111.11188, 2021. URL  
438 <https://arxiv.org/abs/2111.11188>.
- 439 [33] Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamienny, Philip Torr,  
440 Wendelin Böhrer, and Shimon Whiteson. Facmac: Factored multi-agent centralised policy  
441 gradients. *Advances in neural information processing systems*, 34:12208–12221, 2021.
- 442 [34] Julien Perolat, Joel Z. Leibo, Vinicius Zambaldi, Charles Beattie, Karl Tuyls, and Thore Graepel.  
443 A multi-agent reinforcement learning model of common-pool resource appropriation.  
444 In *Proceedings of the 31st International Conference on Neural Information Processing Systems*,  
445 NIPS’17, pp. 3646–3655, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN  
446 9781510860964.
- 447 [35] Dean A. Pomerleau. Efficient training of artificial neural networks for autonomous navigation.  
448 *Neural Computation*, 3(1):88–97, 1991.
- 449 [36] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas  
450 Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon  
451 Whiteson. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- 452 [37] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
453 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 454 [38] Callum Rhys Tilbury, Juan Claude Formanek, Louise Beyers, Jonathan Phillip Shock, and Arnu  
455 Pretorius. Coordination failure in cooperative offline MARL. In *ICML 2024 Workshop: Aligning  
456 Reinforcement Learning Experimentalists and Theorists*, 2024. URL <https://openreview.net/forum?id=gR731hgNB1>.  
457
- 458 [39] Michael Tomasello, Malinda Carpenter, Josep Call, Tanya Behne, and Henrike Moll. Un-  
459 derstanding and sharing intentions: The origins of cultural cognition. *Behavioral and Brain  
460 Sciences*, 28(5):675–691, 2005. DOI: 10.1017/S0140525X05000129.
- 461 [40] Xihuai Wang et al. ZSC-Eval: An evaluation toolkit and benchmark for multi-agent zero-shot  
462 coordination. In *Advances in Neural Information Processing Systems*, 2024.
- 463 [41] Daniël Willemsen, Mario Coppola, and Guido C.H.E. de Croon. Mambpo: Sample-efficient  
464 multi-robot reinforcement learning using learned world models. In *2021 IEEE/RSJ International  
465 Conference on Intelligent Robots and Systems (IROS)*, pp. 5635–5640, 2021. DOI: 10.1109/  
466 IROS51168.2021.9635836.

467 **Appendix contents**

- 468 • Appendix A: Benchmark reference tables
- 469 • Appendix B: Per-scenario return histograms
- 470 • Appendix C: Hyperparameters
- 471 • Appendix D: Per-scenario Benchmark 1 vs Benchmark 2 eval returns
- 472 • Appendix E: Per-scenario Benchmark 3 eval returns by split
- 473 • Appendix F: Per-scenario background-agent return change
- 474 • Appendix G: Cross-substrate aggregate scores

475 **A Benchmark reference tables**

476 This appendix lists, for reference, the five substrates that compose the Moltenpot benchmark, the full  
 477 per-scenario catalogue with its social-role badges, and the Benchmark 3 train/test splits.

Table 2: Moltenpot substrates and the primary social challenge each substrate is intended to probe.

Substrate	Players	Social challenge
Clean Up	7	Public-goods maintenance, contribution vs. free-riding, reciprocity, sanctioning, and corrigibility under collective dependence.
Coins [26]	2	Reciprocity and adaptation to partner policy in a mixed-motive dyadic dilemma.
Coop Mining	6	Trust formation, partner choice.
Commons Harvest [34]	7	Tragedy of the commons.
Allelopathic Harvest	4–16	Convention adoption and norm-setting: no dominant individual strategy, so the payoff-optimal planting colour depends on what the rest of the population is doing.

Table 3: Scenario overview for the Moltenpot benchmark. **V** denotes visitor-mode focal populations and **R** denotes resident-mode focal populations. **PRO** prosocial / sustainable counterpart, and **ANTI** antisocial / unsustainable counterpart. Primary scenario families are: **RC** reciprocity / conditional cooperation, **PC** partner choice / partner discrimination, **CG** commons governance / sustainability, **TR** trust / institutional commitment, and **CA** convention adoption. Modifier badges indicate the dominant complication: **TH** harsher (low) reciprocity threshold / suspiciousness, **EN** punishment / sanctioning, **CO** corrigibility, **NS** non-stationary or alternating partner behavior, **LT** latent partner type, and **ST** structured.

Substrate	Sc.	Scenario Type	Description
Clean Up	0	<b>V</b> <b>PRO</b>	Visiting an altruistic population.
	2	<b>V</b> <b>ST</b>	Visiting a turn-taking population that cleans first.
	3	<b>V</b> <b>ST</b>	Visiting a turn-taking population that eats first.
	4	<b>R</b> <b>RC</b>	Focals are visited by one reciprocator.
	5	<b>R</b> <b>RC</b> <b>TH</b>	Focals are visited by two suspicious reciprocators.
	7	<b>V</b> <b>RC</b> <b>TH</b>	Focals visit a resident group of suspicious reciprocators.
	13	<b>V</b> <b>RC</b> <b>CO</b>	Focals visit easily corrigible reciprocators.
	14	<b>V</b> <b>RC</b> <b>CO</b> <b>TH</b>	Focals visit reciprocators who are corrigible but difficult to convince.
Coins	15	<b>V</b> <b>RC</b> <b>CO</b> <b>NS</b>	Focals visit easily corrigible reciprocators and also bots who alternate between contributing and free riding.
	0	<b>PC</b> <b>LT</b>	Partner is either a pure cooperators or a pure defector.
	1	<b>RC</b>	Partner is a high-threshold (generous) reciprocator.
	2	<b>RC</b> <b>TH</b>	Partner is a low-threshold (harsh) reciprocator.
	3	<b>RC</b> <b>EN</b>	Partner is a high-threshold (generous) strong reciprocator.
	4	<b>RC</b> <b>EN</b> <b>TH</b>	Partner is a low-threshold (harsh) strong reciprocator.
	5	<b>PC</b> <b>PRO</b>	Partner is a cooperators.
6	<b>PC</b> <b>ANTI</b>	Partner is a defector.	
Coop Mining	0	<b>V</b> <b>TR</b> <b>PRO</b>	Visiting cooperators.
	1	<b>V</b> <b>PC</b>	Visiting residents that extract both ores.
	2	<b>V</b> <b>PC</b> <b>ANTI</b>	Visiting defectors.
	3	<b>R</b> <b>PC</b> <b>PRO</b>	Residents visited by a cooperators.
	4	<b>R</b> <b>PC</b> <b>ANTI</b>	Residents visited by a defector.
Commons Harvest: Closed	5	<b>V</b> <b>PC</b> <b>LT</b>	Find the cooperators partner.
	0	<b>V</b> <b>CG</b> <b>ANTI</b>	Focals visit pacifist bots who harvest unsustainably.
	1	<b>R</b> <b>CG</b> <b>ANTI</b>	Focals are resident and visited by pacifist bots who harvest unsustainably.
	2	<b>V</b> <b>CG</b> <b>PRO</b> <b>EN</b>	Focals visit bots who zap and harvest sustainably if they get a chance.
Commons Harvest: Open	3	<b>R</b> <b>CG</b> <b>PRO</b> <b>EN</b>	Focals are resident, and are visited by bots who zap and harvest sustainably if they get a chance.
	0	<b>R</b> <b>CG</b> <b>ANTI</b> <b>EN</b>	Focals are resident and visited by two bots who zap and harvest unsustainably.

Table 3 (continued)

Substrate	Sc.	Scenario Type	Description
	1	<b>R</b> <b>CG</b> <b>ANTI</b>	Focals are resident and visited by two pacifists who harvest unsustainably.
Commons Harvest: Partnership	0	<b>V</b> <b>TR</b> <b>PRO</b>	Meeting good partners. 1 focal agent.
	1	<b>R</b> <b>TR</b> <b>PRO</b>	Focals are resident and visitors are good partners.
	2	<b>V</b> <b>TR</b> <b>PRO</b> <b>EN</b>	Focals visit zappers who harvest sustainably but lack trust.
	3	<b>R</b> <b>TR</b> <b>PRO</b> <b>EN</b>	Focals are resident and visited by zappers who harvest sustainably but lack trust.
	4	<b>V</b> <b>TR</b> <b>ANTI</b>	Focals visit pacifists who do not harvest sustainably.
	5	<b>V</b> <b>TR</b> <b>ANTI</b> <b>EN</b>	Focals visit zappers who do not harvest sustainably.
Allelopathic Harvest	0	<b>V</b> <b>CA</b> <b>ANTI</b>	Visiting a population where planting green berries is the prevailing convention.
	1	<b>V</b> <b>CA</b> <b>PRO</b>	Visiting a population where planting red berries is the prevailing convention.
	2	<b>R</b> <b>CA</b>	Focals are resident and visited by bots who plant either red or green.
	3	<b>CA</b> <b>PRO</b>	Focals like red, visited by convention followers.
	4	<b>CA</b> <b>PRO</b>	Focals like red, visited by mixture of convention followers and bots who like red.
	5	<b>CA</b> <b>ANTI</b>	Focals like red, visited by mixture of convention followers and bots who like green.
	6	<b>CA</b> <b>ANTI</b>	Small-world version of scenario 5 (5 focal, 5 bots).
	7	<b>CA</b> <b>PRO</b>	Small-world pure-follower version of scenario 3 (5 focal, 5 bots).
	8	<b>CA</b> <b>PRO</b>	Minimal 2-focal/2-follower setting: focals like red, meeting convention followers.
	9	<b>R</b> <b>CA</b>	Focal majority with mixed preferences; focals are resident, visited by a couple of convention followers.
	10	<b>R</b> <b>CA</b>	Very small mixed-preference focal group with a single follower bot.
	11	<b>R</b> <b>CA</b>	Large focal majority (8R + 6G) visited by a couple of convention followers.
12	<b>R</b> <b>CA</b> <b>ANTI</b>	Largest focal majority visited by a couple of convention followers and bots who plant green.	

Table 4: Moltenpot Evaluation 3 train–test splits. For each substrate, the split specifies a disjoint ( $\mathcal{S}^{\text{tr}}$ ,  $\mathcal{S}^{\text{te}}$ ) pair and the social shift it is designed to probe.

Substrate	Split Shift	Train	Test	What this split tests
Clean Up	U1 <b>PRO</b> <b>ST</b> → <b>ANTI</b>	0, 2, 3, 17	1, 21, 22	Generalisation from helpful or structured partner populations to free-rider intrusion.
	U2 <b>RC</b> → <b>RC</b> <b>TH</b>	4, 8, 9, 11, 13, 15, 18	5, 6, 7, 14, 16	Transfer from low-threshold or nice reciprocity to suspicious, higher-threshold reciprocity.
	U3 <b>RC</b> → <b>RC</b> <b>CO</b>	4, 5, 6, 7, 8, 18	9, 11, 12, 13, 14, 15, 16	Generalisation from reciprocity to corrigibility under partner variation.
	U4 <b>PRO</b> <b>RC</b> <b>ST</b> → <b>NS</b> <b>EN</b>	0, 2, 3, 4, 5, 6, 7, 8, 17, 18	10, 15, 16, 19, 20	Robustness to unstable and sanctioning partner behaviour after training on stable helpful and reciprocal counterparts.
Coins	C1 <b>PC</b> → <b>RC</b>	0, 5, 6	1, 2, 3, 4	Generalisation from unconditional partner types to conditional reciprocity.
	C2 <b>RC</b> → <b>RC</b> <b>TH</b>	1, 3	2, 4	Transfer from lenient to harsh reciprocity thresholds.
	C3 <b>RC</b> → <b>RC</b> <b>EN</b>	1, 2	3, 4	Transfer across punishment severity.
	C4 <b>PC</b> → <b>PC</b> <b>LT</b>	5, 6	0	Latent partner-type inference under within-episode uncertainty.
Coop Mining	M1 <b>V</b> → <b>R</b>	0, 1, 2, 5	3, 4	Transfer from visitor-mode response to resident-mode competition.
	M2 <b>PC</b> → <b>PC</b> <b>LT</b>	0, 2	5	Generalisation from homogeneous partner populations to heterogeneous populations.
	M3 <b>V</b> <b>PC</b> → <b>R</b> <b>PC</b>	0, 5	3	Generalisation from abundant cooperation opportunities to scarcity and competition for partner access.
Commons Harvest	H1 <b>ANTI</b> → <b>ANTI</b> <b>EN</b>	Closed: 0, 1; Open: 1; Partnership: 4	Open: 0; Partnership: 5	Transfer from passive unsustainable counterparts to enforcement-capable unsustainable counterparts.
	H2 <b>V</b> → <b>R</b>	Closed: 0, 2; Partnership: 0, 2, 4, 5	Closed: 1, 3; Open: 0, 1; Partnership: 1, 3	Transfer from visitor-mode adaptation to resident-mode governance.
	H3 <b>CG</b> → <b>TR</b>	Closed: 0, 1, 2, 3; Open: 0, 1	Partnership: 0, 1, 2, 3	Generalisation from anonymous commons management to dyadic institutional commitment.
	H4 <b>TR</b> <b>PRO</b> → <b>TR</b> <b>PRO</b> <b>EN</b>	Partnership: 0, 1	Partnership: 2, 3	Trust calibration under partner shift from good partners to guarded or adversarial ones.
Allelopathic Harvest	A1 <b>PRO</b> → <b>ANTI</b>	1, 3, 4, 7, 8	0, 5, 6	Generalisation from aligned or shapeable conventions to misaligned or opposing conventions.
	A2 pure followers → mixed visitors	3, 7, 8, 9, 10, 11	4, 5, 6, 12	Robustness of convention-following when some visitors no longer merely follow the prevailing norm but also carry their own colour preference, in both visitor-like and resident-majority settings.
	A3 large → small	3, 4, 5, 11, 12	6, 7, 8, 9, 10	Transfer from large-population convention dynamics to sparse settings where a few agents or even a single plant can determine which convention takes hold.
	A4 <b>V</b> → <b>R</b>	0, 1	2, 9, 10, 11, 12	Transfer from adapting to an already-established visiting convention to sustaining or shaping the convention as a resident focal majority under different visitor mixes.

478 **B Per-scenario return histograms**

479 Figures 7–11 show, for every scenario in the benchmark, the distribution of focal-agent episode  
 480 returns across the 1000 trajectories logged during the scenario’s independent-PPO run. Because  
 481 trajectories are recorded uniformly across training time, each histogram should exhibit a wide support  
 482 with a left tail from near-random initial rollouts and a right peak near the converged PPO return;  
 483 uniform sampling across training time deliberately produces a skill-mixed dataset rather than an  
 484 expert-only one, following the dataset-transparency recommendation of Formanek et al. [12]. Narrow  
 485 histograms are a warning sign that the scenario’s behaviour policy collapsed or that the logging  
 486 schedule failed to capture the full training arc. In every panel the dashed red line marks the scenario’s  
 487 mean episode return and the dashed green line marks the 90th percentile.

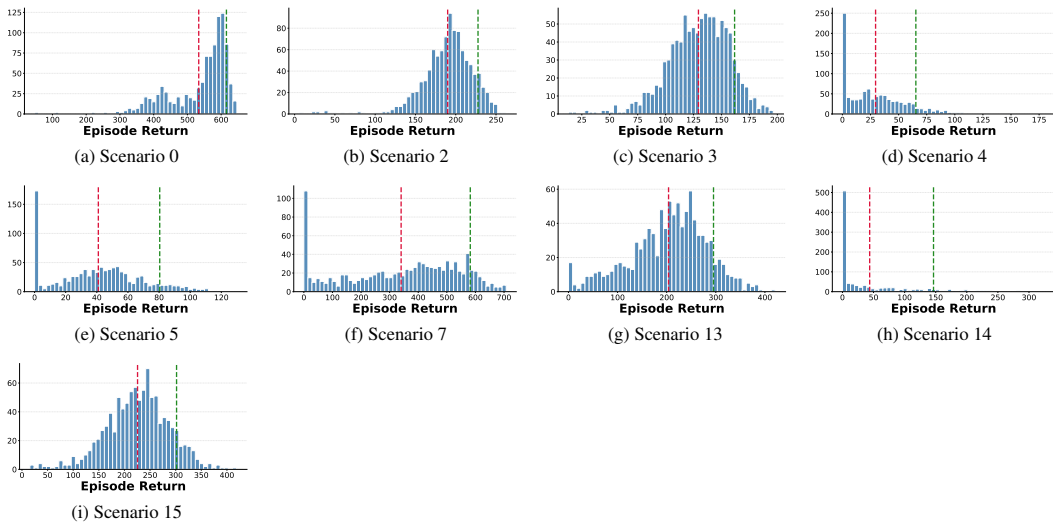


Figure 7: **Clean Up** — per-scenario focal-return distributions. Dashed red: scenario mean. Dashed green: 90th percentile.

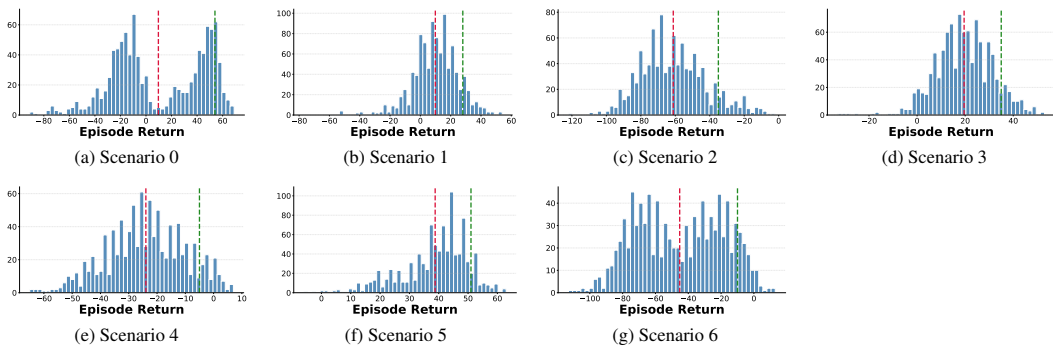


Figure 8: **Coins** — per-scenario focal-return distributions. Dashed red: scenario mean. Dashed green: 90th percentile.

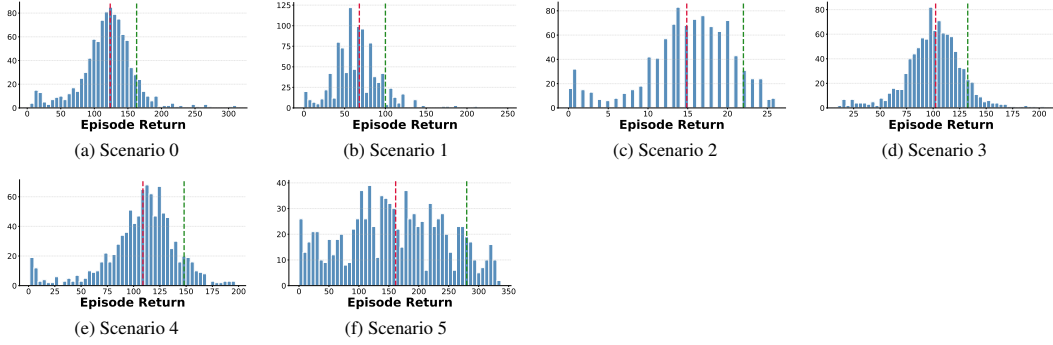


Figure 9: **Coop Mining** — per-scenario focal-return distributions. Dashed red: scenario mean. Dashed green: 90th percentile.

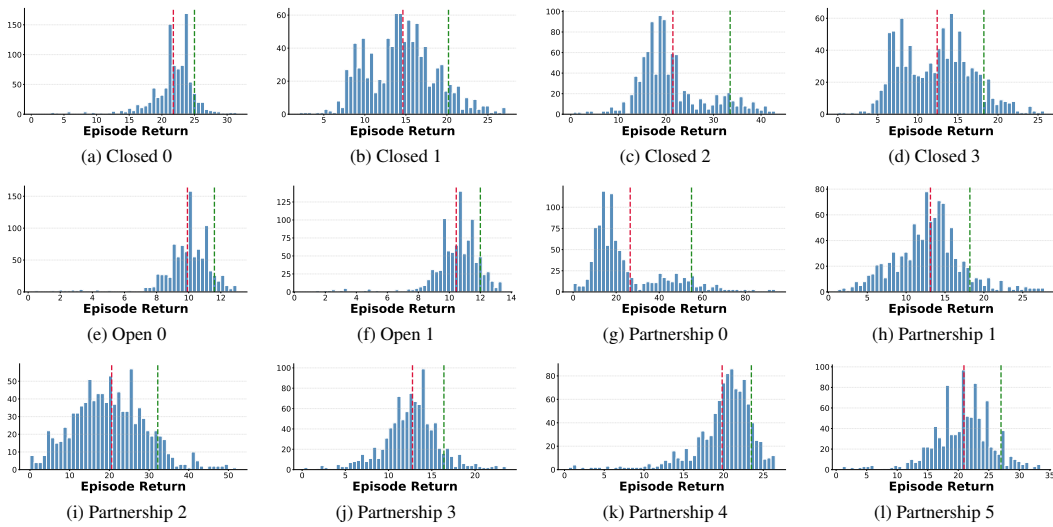


Figure 10: **Commons Harvest** — per-scenario focal-return distributions. Variants are grouped by row: *closed* (top), *open + partnership* (middle and bottom). Dashed red: scenario mean. Dashed green: 90th percentile.

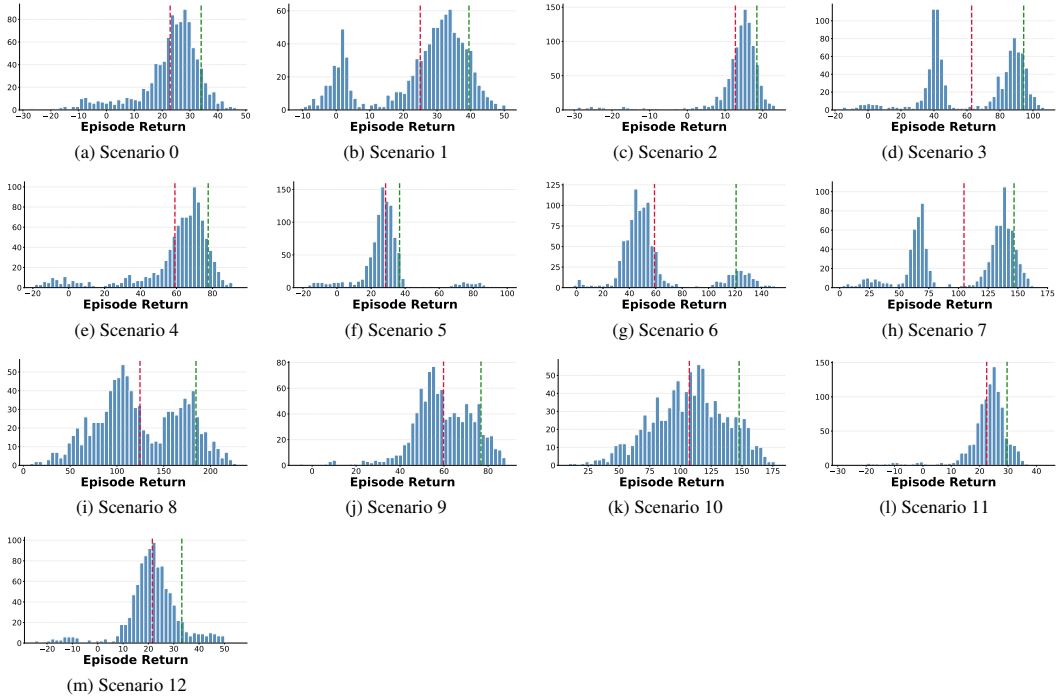


Figure 11: Allelopathic Harvest — per-scenario focal-return distributions. Dashed red: scenario mean. Dashed green: 90th percentile.

488 **C Hyperparameters**

489 **Tuning protocol.** We purposely use a modest hyperparameter tuning budget. For every (algorithm,  
 490 substrate) pair we In particular, hyperparameters are tuned on the *first scenario* of the substrate only,  
 491 and the resulting values are then reused unchanged on every other scenario in the substrate. It is  
 492 now well-recognised that for offline RL the online tuning budget should be kept small and fixed,  
 493 so that comparisons reward algorithms that need less tuning rather than the practitioner who could  
 494 afford the most online evaluations [13]; in any realistic deployment of an offline-trained policy, online  
 495 evaluation for hyperparameter selection is exactly the resource we cannot liberally spend.

496 **Shared settings.** Every algorithm uses the same focal-agent architecture, optimiser and training  
 497 schedule (Table 5). The only substrate-dependent hyperparameters are given in Table 6.

Table 5: Shared training and architecture settings, applied across all four offline algorithms and all five substrates.

Setting	Value
Convolutional backbone	3-layer CNN (matches MeltingPot 2.4)
Hidden size (FC)	256
Recurrent core (GRU)	256
Optimiser	Adam (default $\beta_1, \beta_2$ )
Learning rate	$3 \times 10^{-4}$
Discount $\gamma$	0.99
Soft-target rate $\tau$	0.005
Batch size	32
Sequence length	128
Gradient steps	20,000
Seeds per (substrate, algorithm)	3

Table 6: Per-algorithm, per-substrate hyperparameters. Values were selected by the per-substrate tuning protocol described above (tune on the substrate’s first scenario, then re-use for every other scenario in the substrate).

Algorithm	Hyperparameter	Coins	Coop Mining	Clean Up	Commons Harvest	Allelopathic Harvest
IQL	Expectile $\tau_{IQL}$	0.80	0.80	0.70	0.70	0.70
	Advantage temp. $\beta$	2.0	4.0	5.0	2.0	3.0
BCQ	Action threshold $\tau_{BCQ}$	0.30	0.30	0.30	0.30	0.40
CQL	CQL coefficient $\alpha$	3.0	5.0	7.0	3.0	6.0

498 **D Per-scenario Benchmark 1 vs Benchmark 2 eval returns**

499 For each scenario in the benchmark we report the side-by-side Benchmark 1 vs Benchmark 2 eval-  
 500 return comparison in raw (un-normalised) return units, organised here as one mosaic per substrate.  
 501 Within each panel the three offline-RL algorithms appear as paired markers: Benchmark 1 as a filled  
 502 circle (single-scenario training, evaluated on that same scenario) and Benchmark 2 as a filled square  
 503 (all-scenarios training, evaluated on that scenario). The marker is the mean across the seeds with data  
 504 and the error bar is the standard error of that mean. Reference lines mark the dataset’s mean episode  
 505 return (red dashed), its  $p_{90}$  “best-in-dataset” ceiling (green dashed), and the random-policy baseline  
 506 (grey dotted). Y-axes are in raw return units, scaled per-panel so within-scenario differences between  
 507 algorithms remain legible across scenarios with very different return magnitudes.

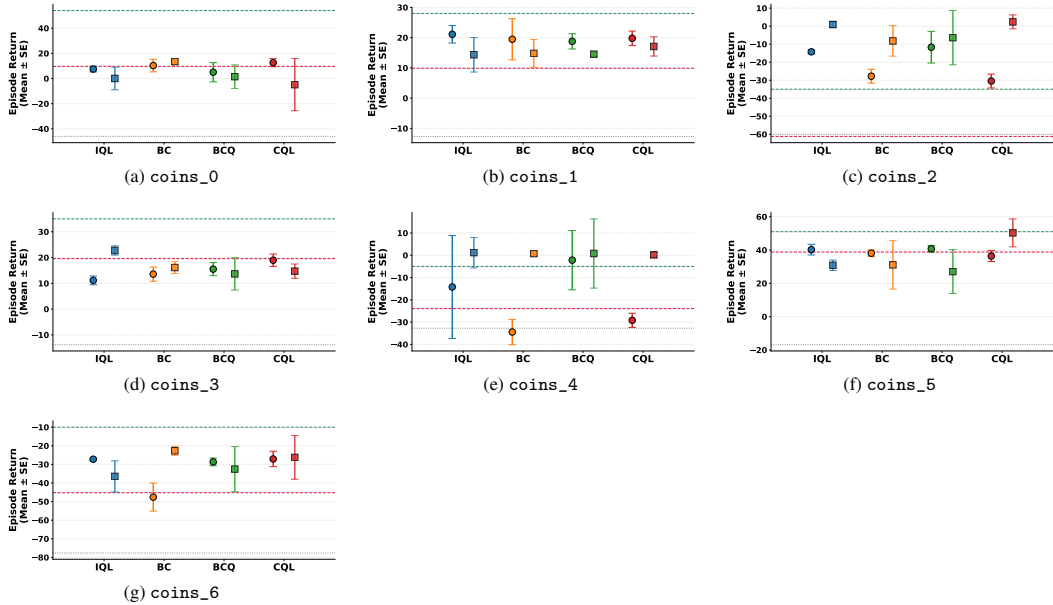


Figure 12: **Per-scenario Benchmark 1 vs Benchmark 2 raw eval returns — Coins.** Each panel shows one scenario in Coins; within each panel the three offline-RL algorithms appear as paired markers — Benchmark 1 (filled circle, single-scenario training) and Benchmark 2 (filled square, all-scenarios training, slightly transparent). Marker = mean across the seeds with data, error bar = SE. Reference lines mark the dataset’s mean episode return (red dashed),  $p_{90}$  “best-in-dataset” ceiling (green dashed), and the random-policy baseline (grey dotted). Y-axes are in raw return units (no normalisation), scaled per-panel.

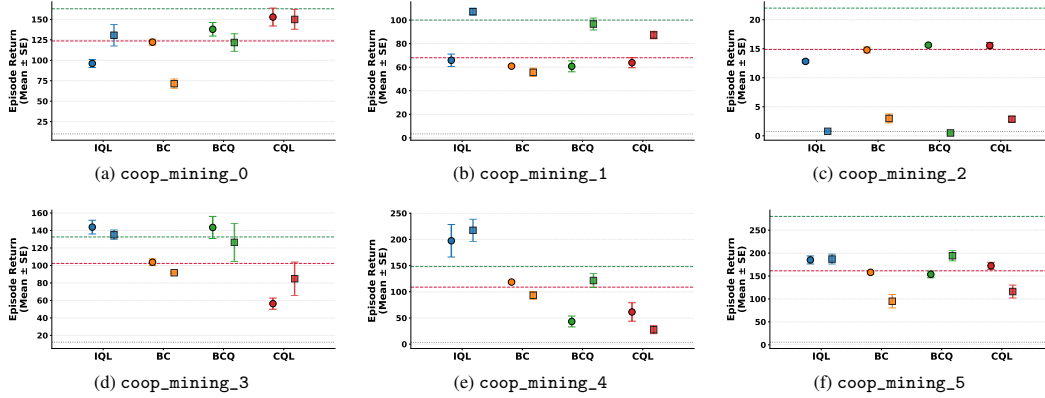


Figure 13: **Per-scenario Benchmark 1 vs Benchmark 2 raw eval returns — Coop Mining.** Each panel shows one scenario in Coop Mining; within each panel the three offline-RL algorithms appear as paired markers — Benchmark 1 (filled circle, single-scenario training) and Benchmark 2 (filled square, all-scenarios training, slightly transparent). Marker = mean across the seeds with data, error bar = SE. Reference lines mark the dataset’s mean episode return (red dashed),  $p_{90}$  “best-in-dataset” ceiling (green dashed), and the random-policy baseline (grey dotted). Y-axes are in raw return units (no normalisation), scaled per-panel.

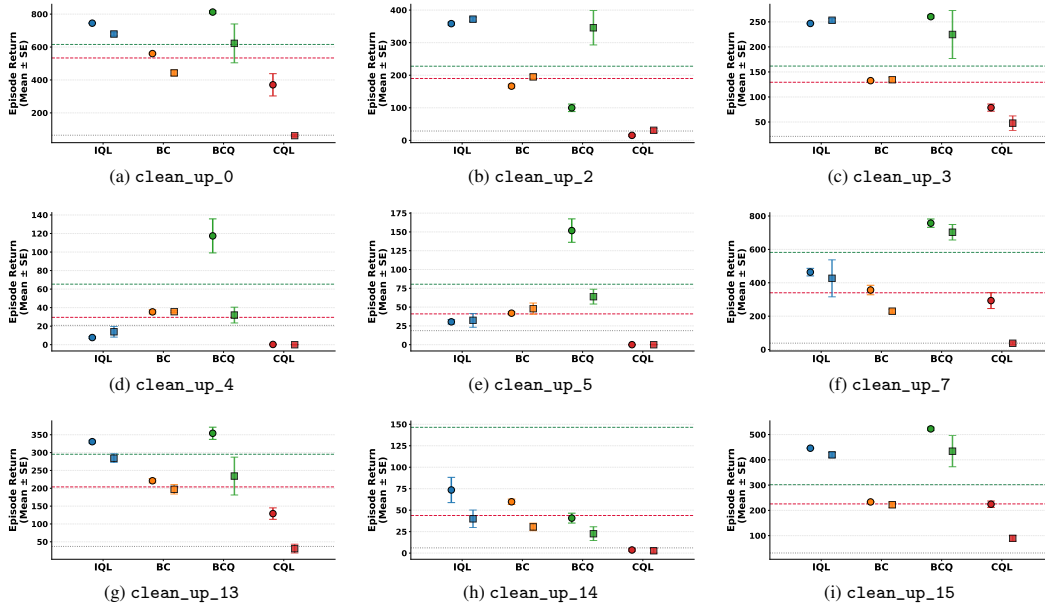


Figure 14: **Per-scenario Benchmark 1 vs Benchmark 2 raw eval returns — Clean Up.** Each panel shows one scenario in Clean Up; within each panel the three offline-RL algorithms appear as paired markers — Benchmark 1 (filled circle, single-scenario training) and Benchmark 2 (filled square, all-scenarios training, slightly transparent). Marker = mean across the seeds with data, error bar = SE. Reference lines mark the dataset’s mean episode return (red dashed),  $p_{90}$  “best-in-dataset” ceiling (green dashed), and the random-policy baseline (grey dotted). Y-axes are in raw return units (no normalisation), scaled per-panel.

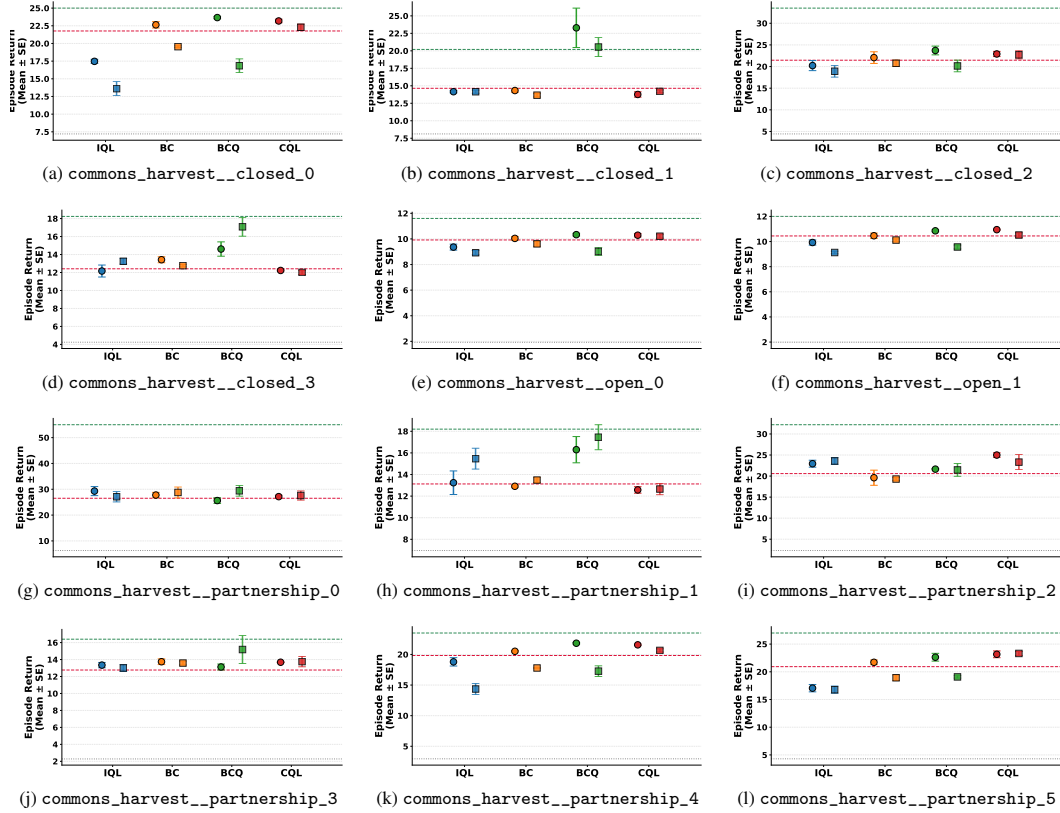


Figure 15: **Per-scenario Benchmark 1 vs Benchmark 2 raw eval returns — Commons Harvest.** Each panel shows one scenario in Commons Harvest; within each panel the three offline-RL algorithms appear as paired markers — Benchmark 1 (filled circle, single-scenario training) and Benchmark 2 (filled square, all-scenarios training, slightly transparent). Marker = mean across the seeds with data, error bar = SE. Reference lines mark the dataset’s mean episode return (red dashed),  $p_{90}$  “best-in-dataset” ceiling (green dashed), and the random-policy baseline (grey dotted). Y-axes are in raw return units (no normalisation), scaled per-panel.

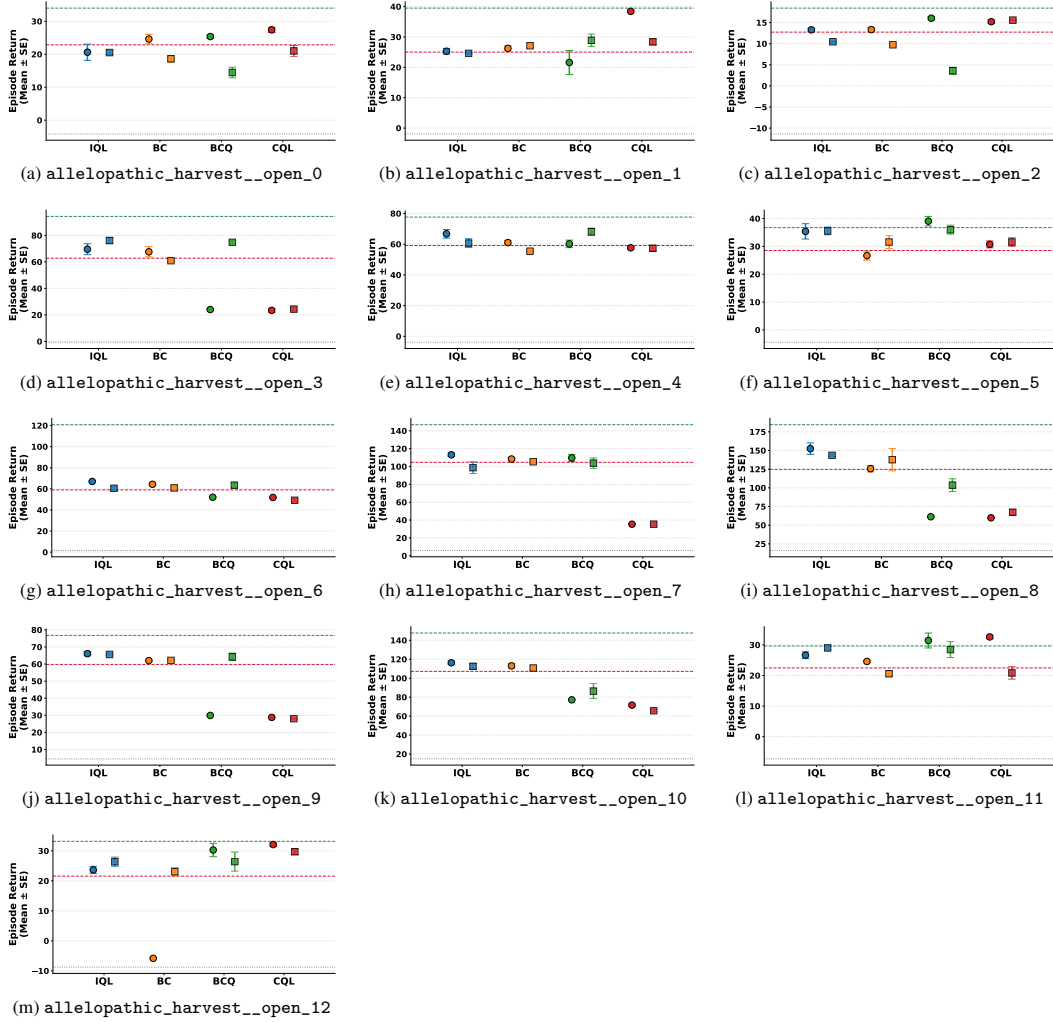
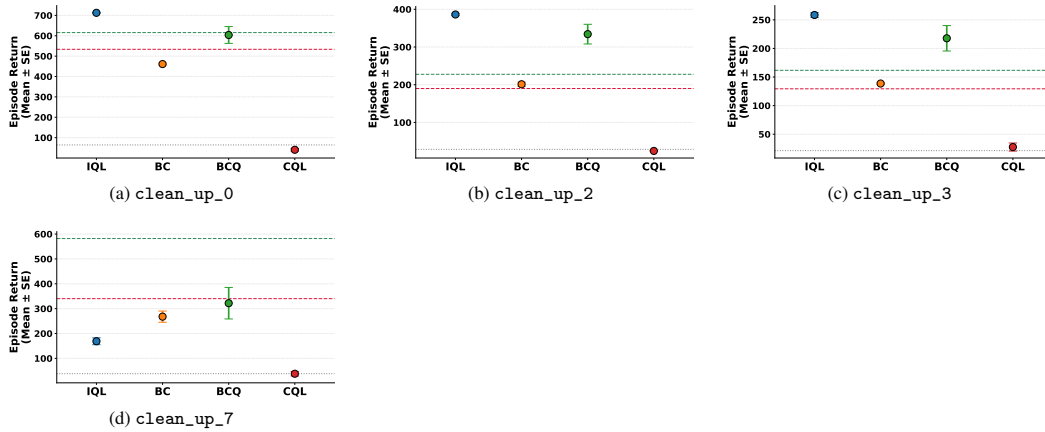


Figure 16: **Per-scenario Benchmark 1 vs Benchmark 2 raw eval returns — Allelopathic Harvest.** Each panel shows one scenario in Allelopathic Harvest; within each panel the three offline-RL algorithms appear as paired markers — Benchmark 1 (filled circle, single-scenario training) and Benchmark 2 (filled square, all-scenarios training, slightly transparent). Marker = mean across the seeds with data, error bar = SE. Reference lines mark the dataset’s mean episode return (red dashed),  $p_{90}$  “best-in-dataset” ceiling (green dashed), and the random-policy baseline (grey dotted). Y-axes are in raw return units (no normalisation), scaled per-panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

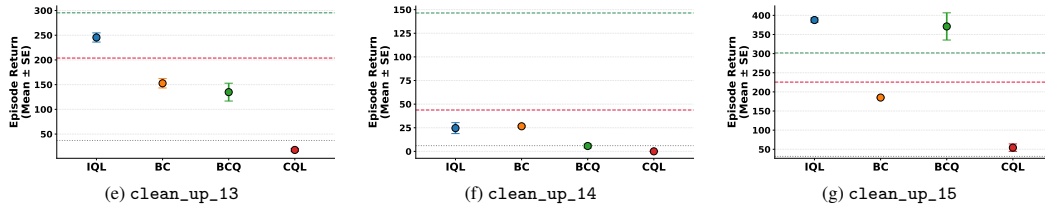
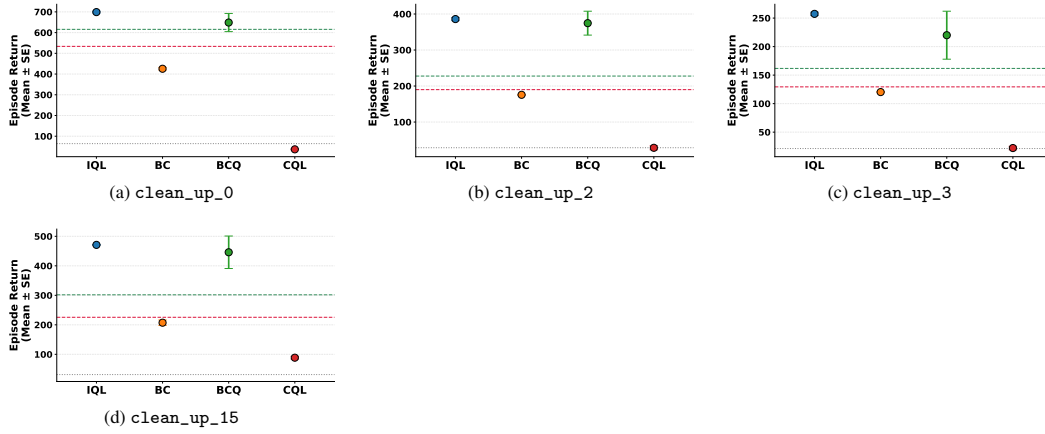


Figure 17: **Split U1 (Clean Up) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

## 508 E Per-scenario Benchmark 3 eval returns by split

509 For each Benchmark 3 train/test split we report the per-scenario final-eval returns in raw (un-  
 510 normalised) return units, with one figure per split. Within each split’s figure, in-distribution (train)  
 511 scenarios are shown above the out-of-distribution (test) scenarios, so the in-/out-of-distribution  
 512 generalisation gap is visible at a glance. Each panel shows the offline-RL algorithms as filled circles,  
 513 where the marker is the mean across the seeds with data and the error bar is the standard error of that  
 514 mean. Reference lines mark the dataset’s mean episode return (red dashed), its  $p_{90}$  “best-in-dataset”  
 515 ceiling (green dashed), and the random-policy baseline (grey dotted). Y-axes are in raw return units,  
 516 scaled per-panel so within-scenario differences between algorithms remain legible across scenarios  
 517 with very different return magnitudes.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

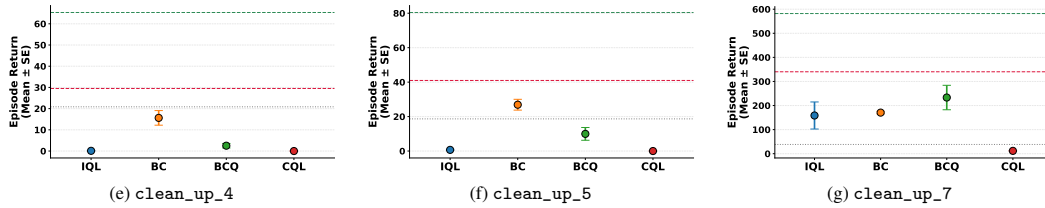
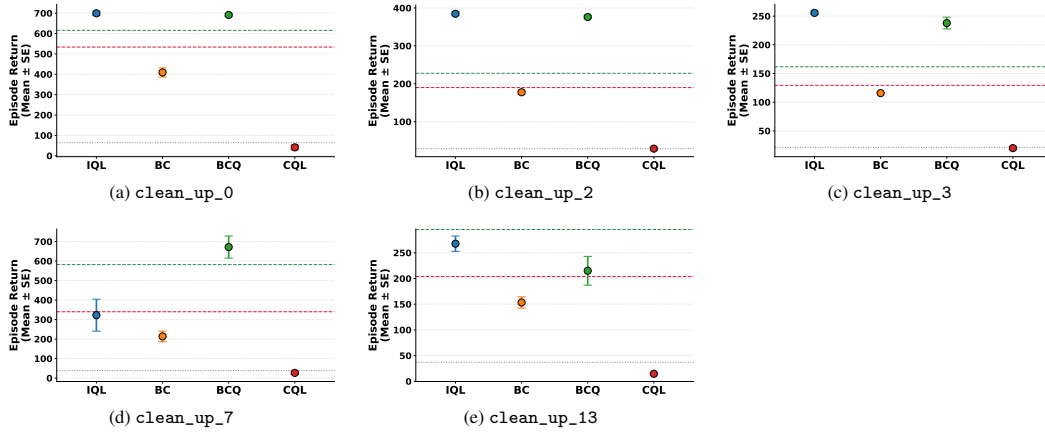


Figure 18: **Split U2 (Clean Up) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

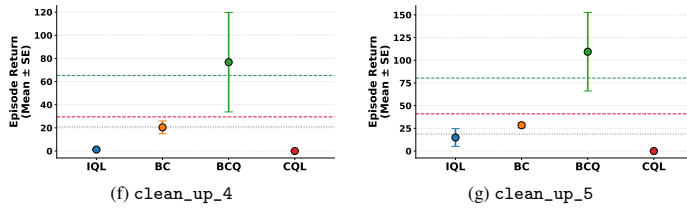
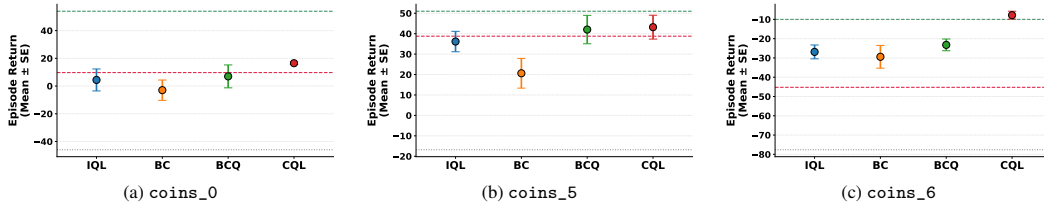


Figure 19: **Split U3 (Clean Up) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

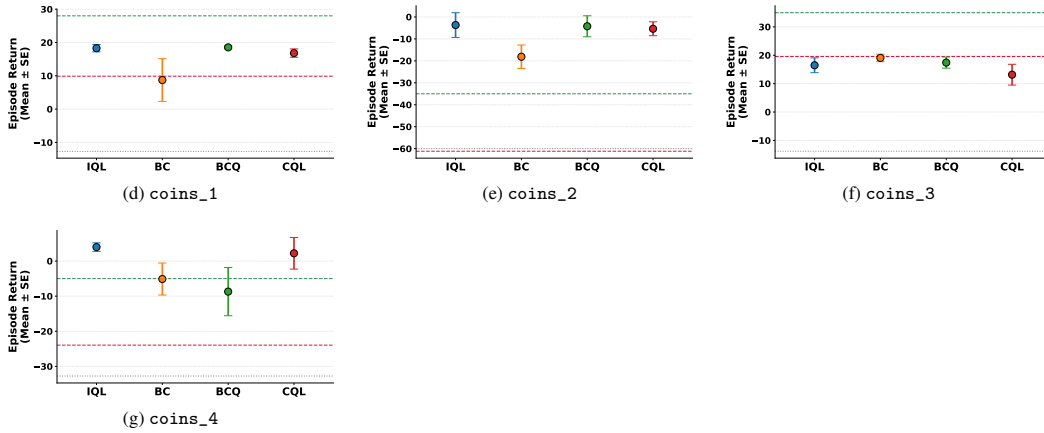
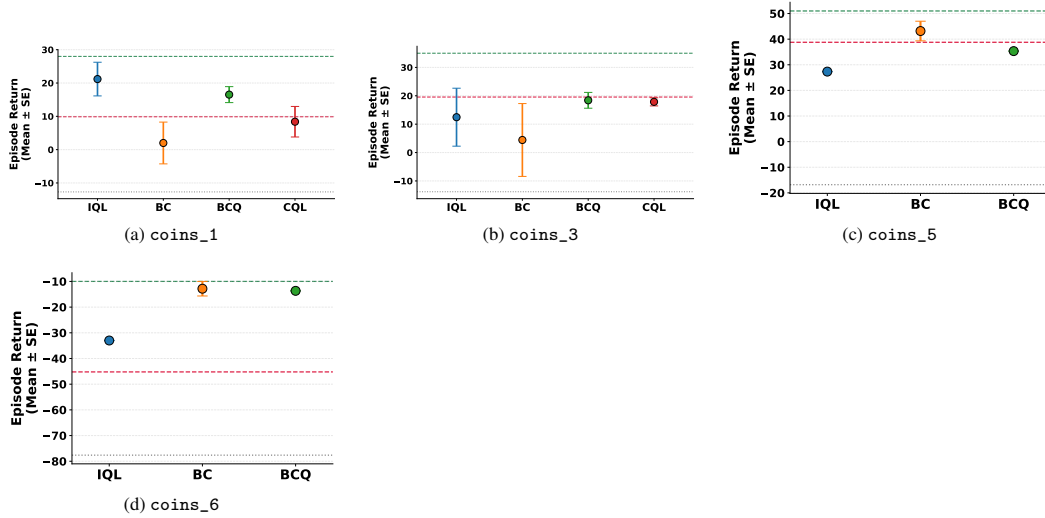


Figure 20: **Split C1 (Coins)** — per-scenario tier-3 eval returns. Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

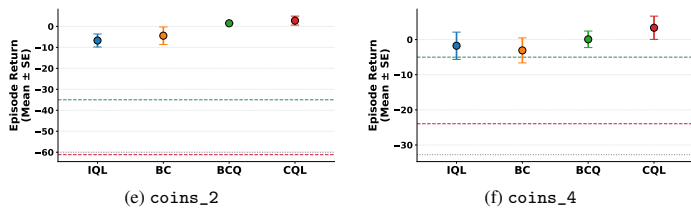
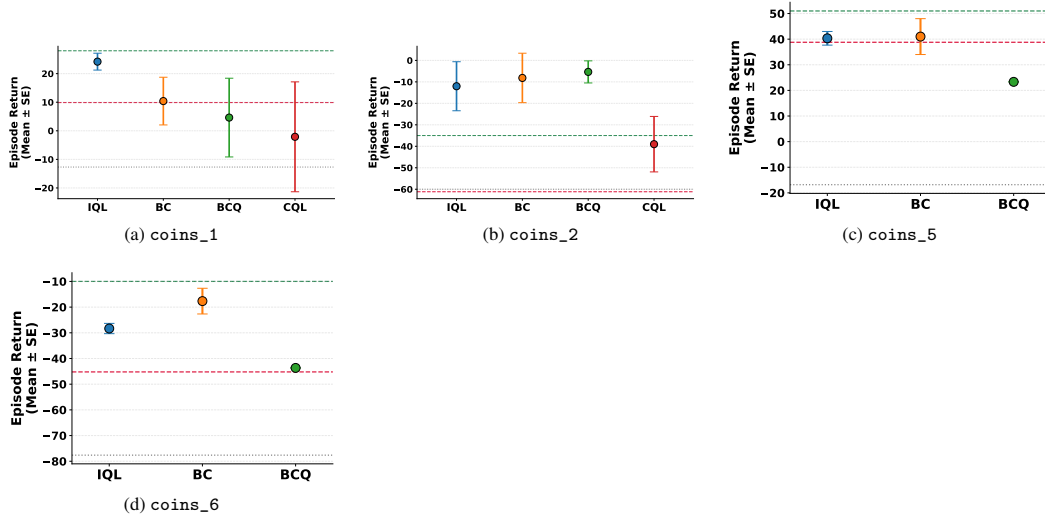


Figure 21: **Split C2 (Coins) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

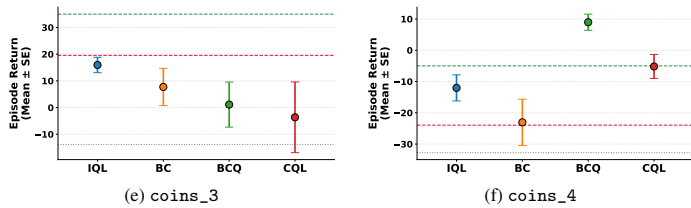
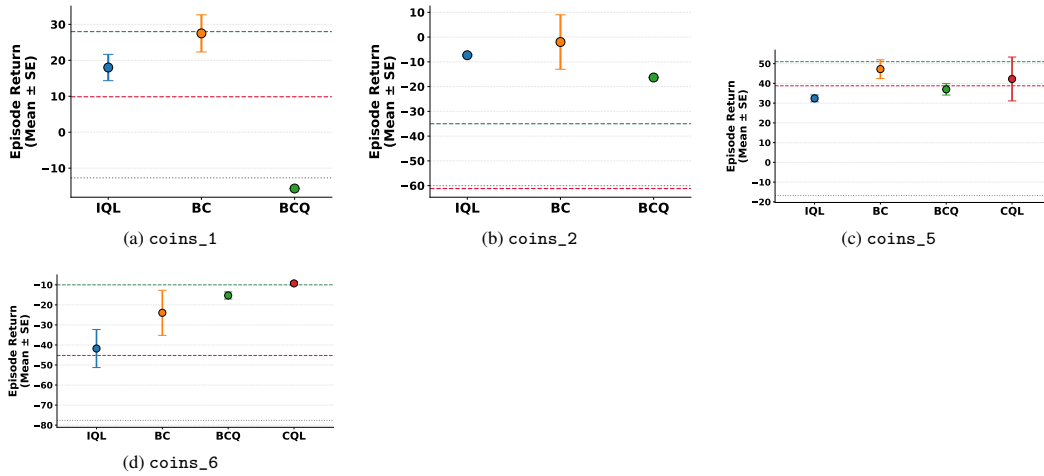


Figure 22: **Split C3 (Coins)** — per-scenario tier-3 eval returns. Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

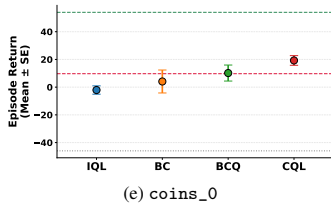
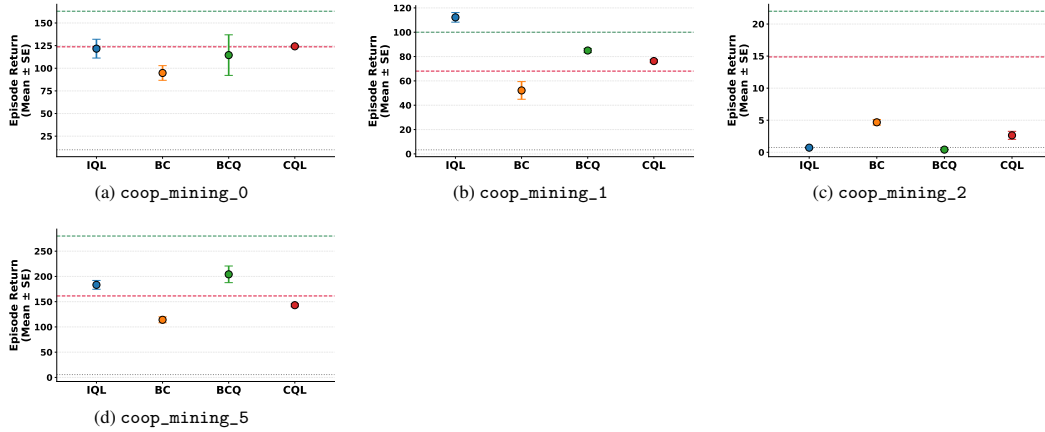


Figure 23: **Split C4 (Coins)** — per-scenario tier-3 eval returns. Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

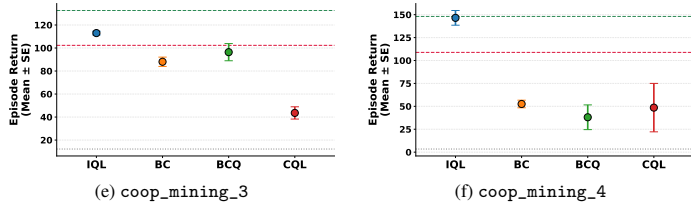
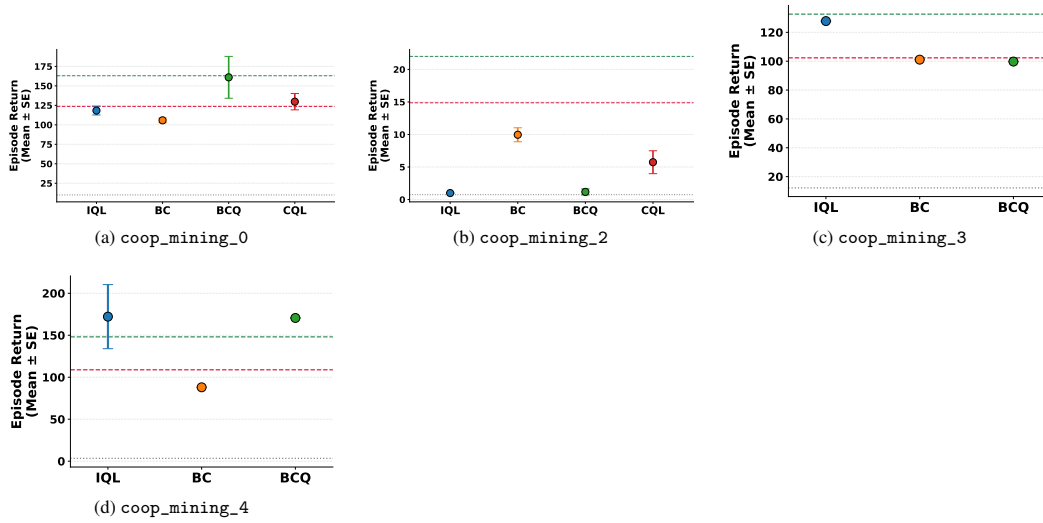


Figure 24: **Split M1 (Coop Mining) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

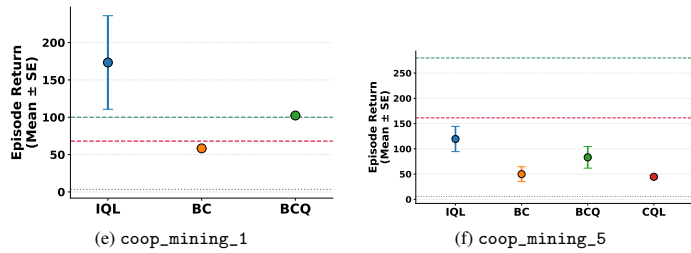
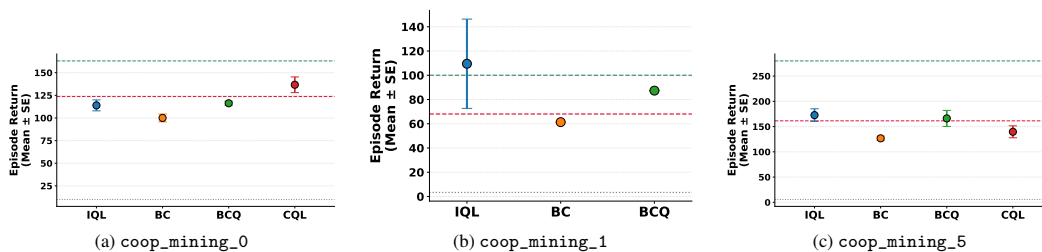


Figure 25: **Split M2 (Coop Mining) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

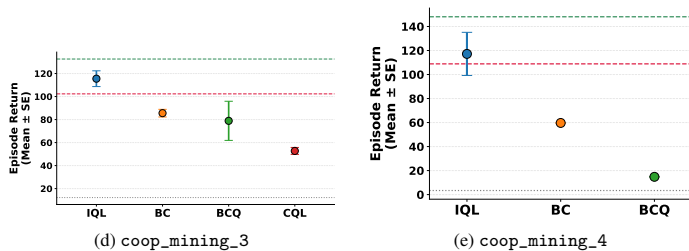
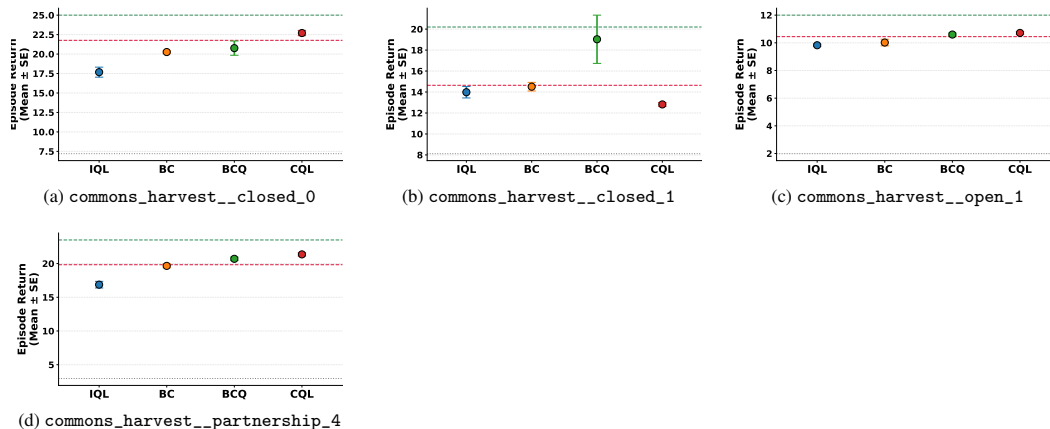


Figure 26: **Split M3 (Coop Mining)** — per-scenario tier-3 eval returns. Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

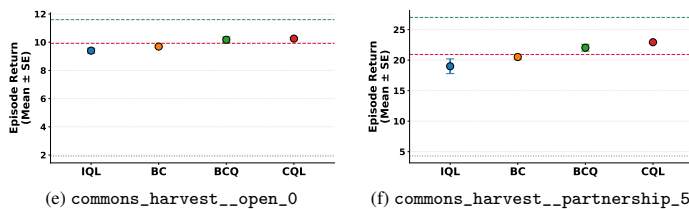
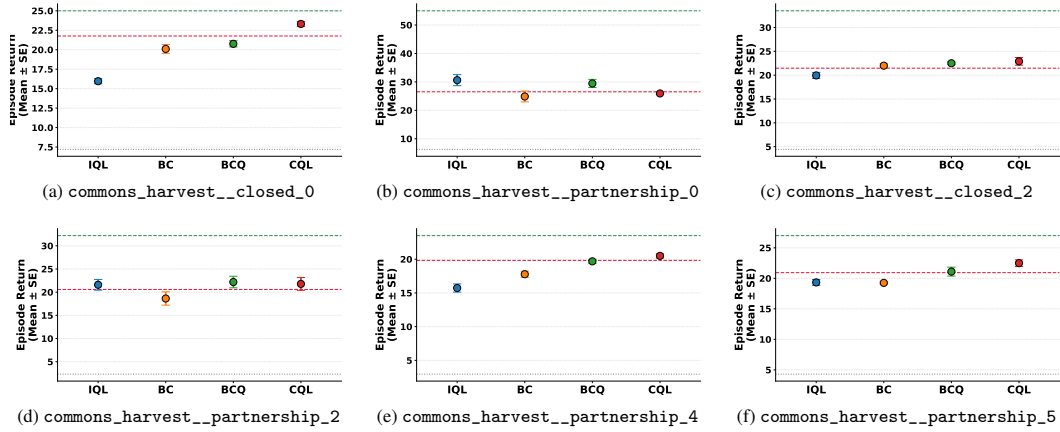


Figure 27: **Split H1 (Commons Harvest)** — per-scenario tier-3 eval returns. Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

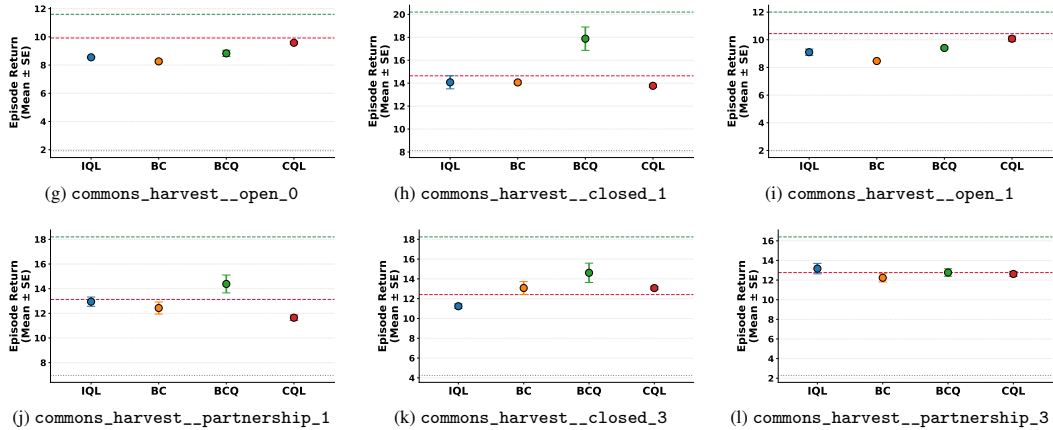
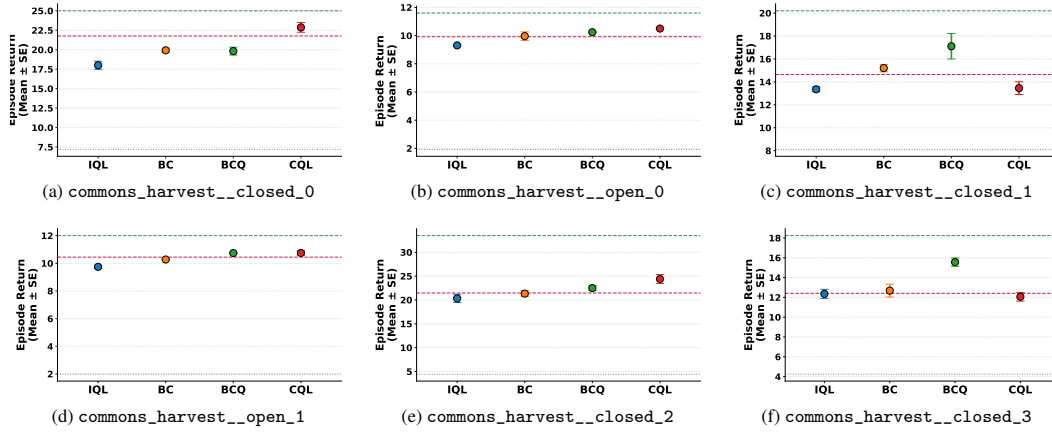


Figure 28: **Split H2 (Commons Harvest) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

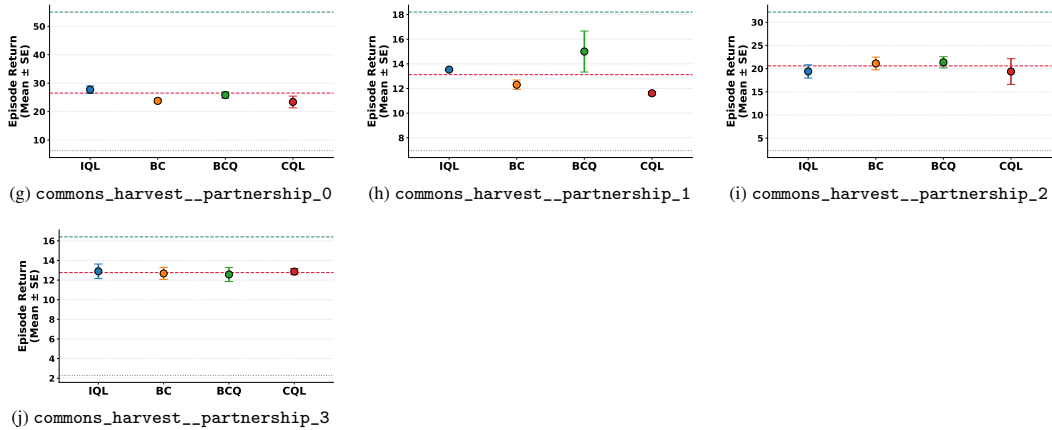
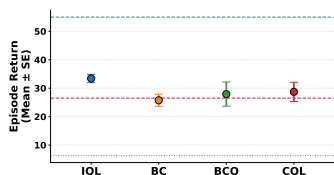
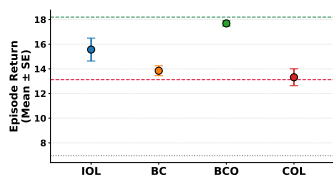


Figure 29: **Split H3 (Commons Harvest)** — per-scenario tier-3 eval returns. Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios

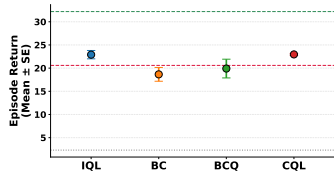


(a) commons\_harvest\_\_partnership\_0

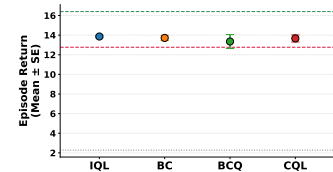


(b) commons\_harvest\_\_partnership\_1

### Out-of-distribution (test) scenarios



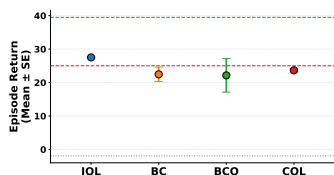
(c) commons\_harvest\_\_partnership\_2



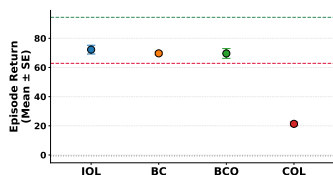
(d) commons\_harvest\_\_partnership\_3

Figure 30: **Split H4 (Commons Harvest) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

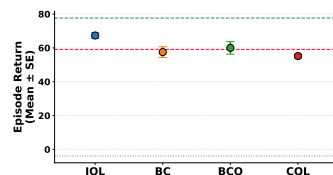
### In-distribution (train) scenarios



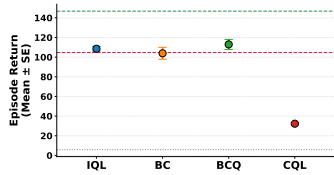
(a) allelopathic\_harvest\_\_open\_1



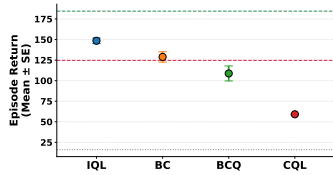
(b) allelopathic\_harvest\_\_open\_3



(c) allelopathic\_harvest\_\_open\_4

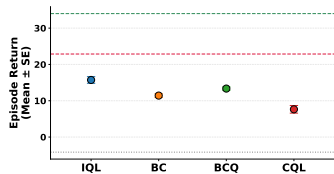


(d) allelopathic\_harvest\_\_open\_7

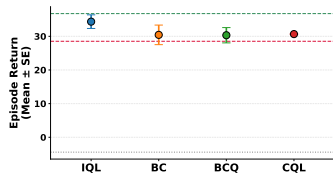


(e) allelopathic\_harvest\_\_open\_8

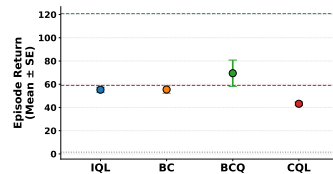
### Out-of-distribution (test) scenarios



(f) allelopathic\_harvest\_\_open\_0



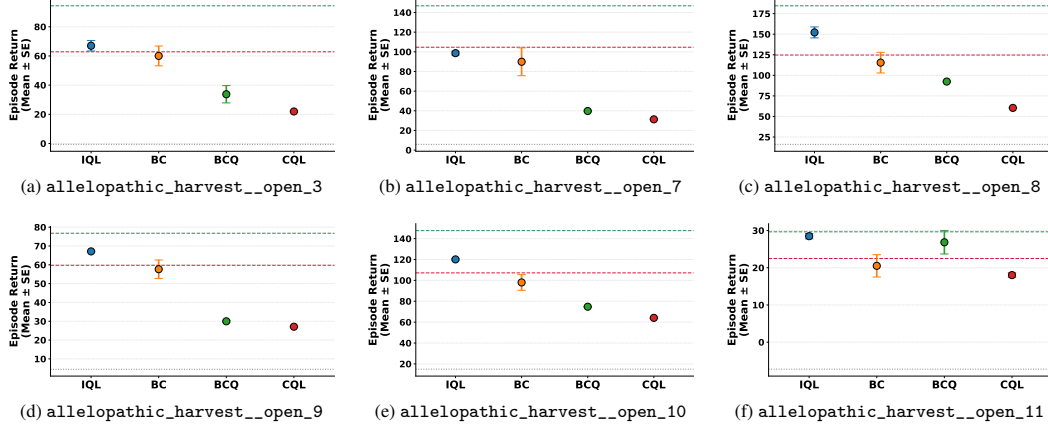
(g) allelopathic\_harvest\_\_open\_5



(h) allelopathic\_harvest\_\_open\_6

Figure 31: **Split A1 (Allelopathic Harvest) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

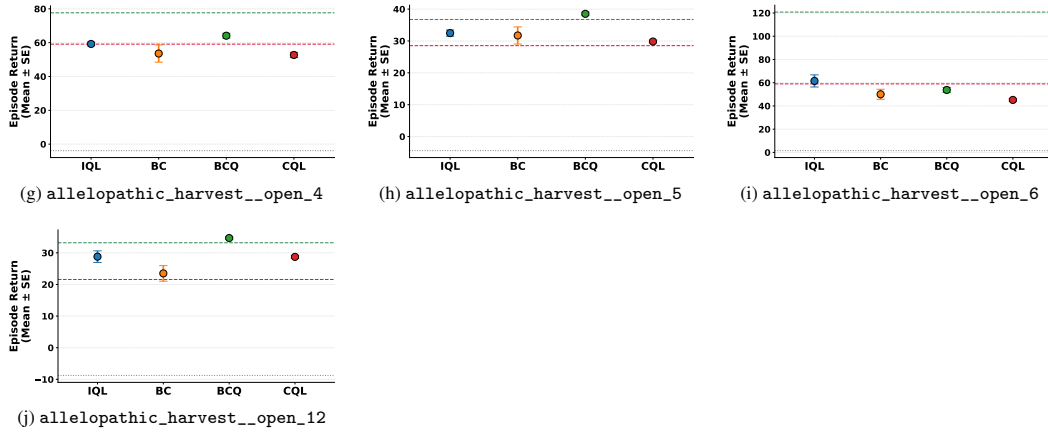
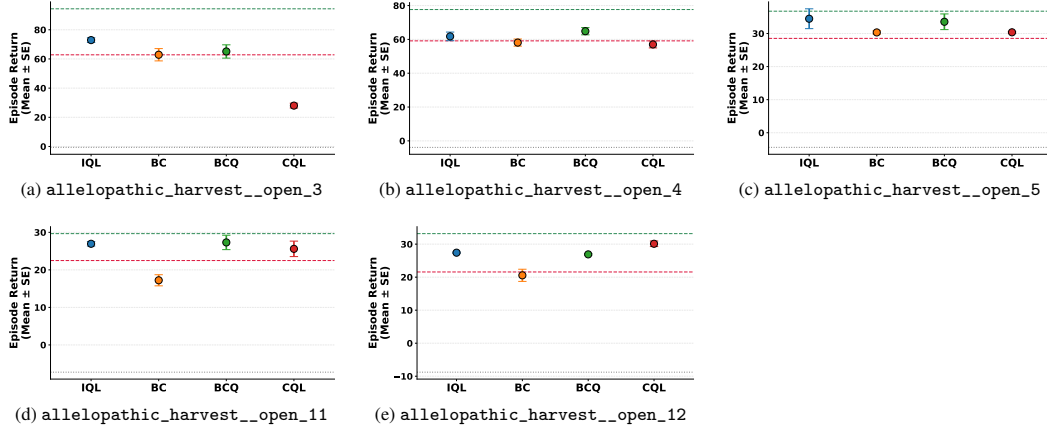


Figure 32: **Split A2 (Allelopathic Harvest)** — per-scenario tier-3 eval returns. Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios



### Out-of-distribution (test) scenarios

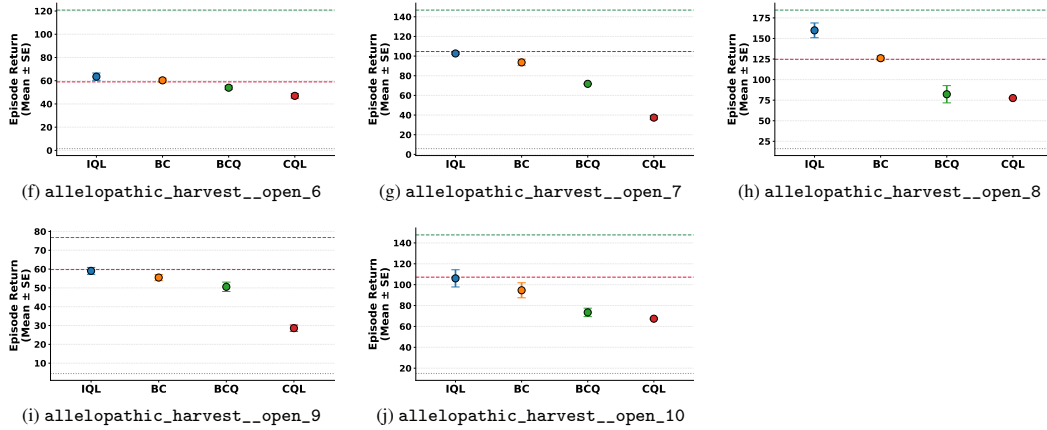
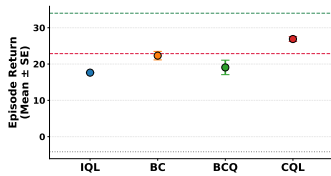
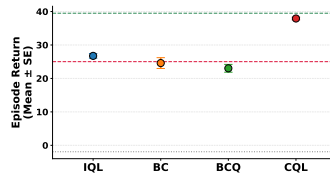


Figure 33: **Split A3 (Allelopathic Harvest)** — per-scenario tier-3 eval returns. Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

### In-distribution (train) scenarios

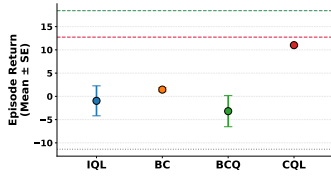


(a) allelopathic\_harvest\_\_open\_0

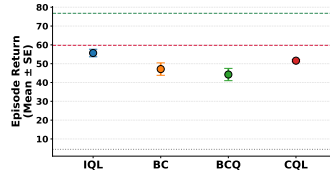


(b) allelopathic\_harvest\_\_open\_1

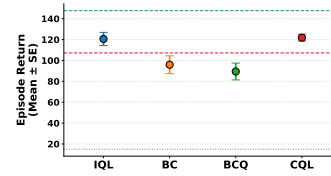
### Out-of-distribution (test) scenarios



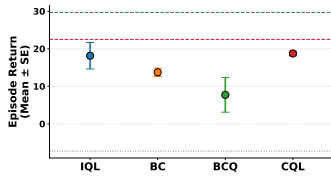
(c) allelopathic\_harvest\_\_open\_2



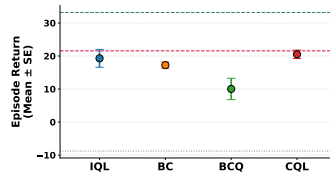
(d) allelopathic\_harvest\_\_open\_9



(e) allelopathic\_harvest\_\_open\_10



(f) allelopathic\_harvest\_\_open\_11



(g) allelopathic\_harvest\_\_open\_12

Figure 34: **Split A4 (Allelopathic Harvest) — per-scenario tier-3 eval returns.** Each panel shows one scenario’s final eval return per algorithm (marker = mean across seeds, error bar = SE). Reference lines: dataset mean (red dashed),  $p_{90}$  best-in-dataset ceiling (green dashed), random-policy baseline (grey dotted). Y-axes are scaled per panel.

518 **F Per-scenario background-agent return change**

519 For each Evaluation 1 scenario we report the per-algorithm change in mean background-agent  
520 episode return induced by the trained focal policy, relative to a uniform-random focal policy  
521 ( $\Delta = R_{\text{trained}} - R_{\text{random}}$ ). Background returns are read from the in-training eval summary  
522 (`eval/in_dist/<scen>/background_return_mean`) of the matching tier-1 source run. Neg-  
523 ative values indicate the focal policy hurt the background population (exploitation); positive values  
524 indicate the focal policy helped them. Panels are grouped by substrate.

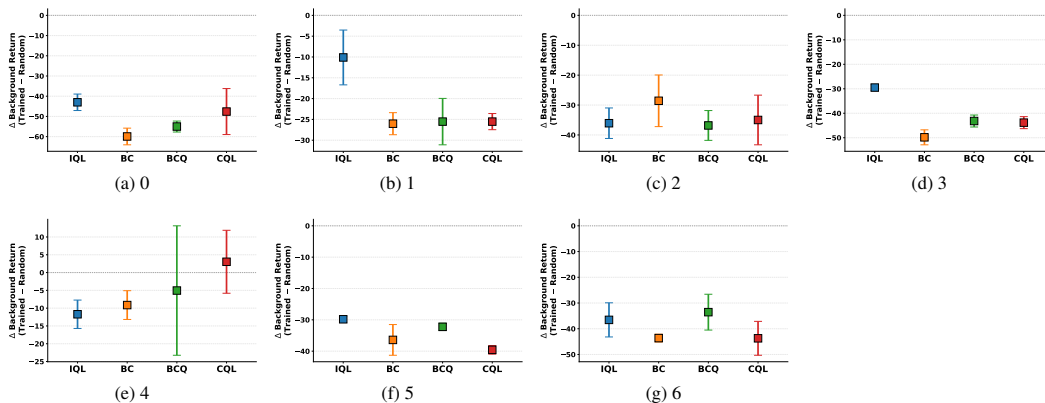


Figure 35: **Coins — background-agent return change.** For each scenario in Coins, one marker per algorithm showing the change in mean background-agent episode return induced by the trained focal policy, relative to a uniform-random focal policy ( $\Delta = R_{\text{trained}} - R_{\text{random}}$ ). Background returns are read from the in-training eval logs. Negative values indicate the focal policy hurt the background population (exploitation); positive values indicate it helped them.

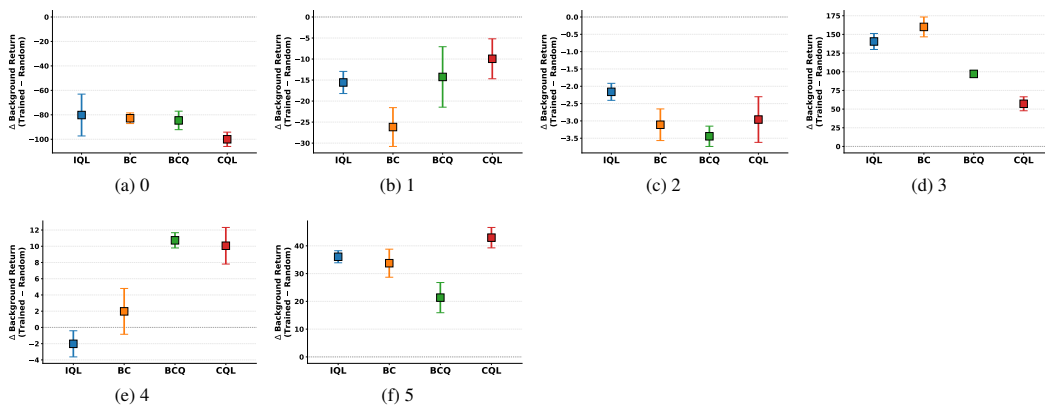


Figure 36: **Coop Mining — background-agent return change.** For each scenario in Coop Mining, one marker per algorithm showing the change in mean background-agent episode return induced by the trained focal policy, relative to a uniform-random focal policy ( $\Delta = R_{\text{trained}} - R_{\text{random}}$ ). Background returns are read from the in-training eval logs. Negative values indicate the focal policy hurt the background population (exploitation); positive values indicate it helped them.

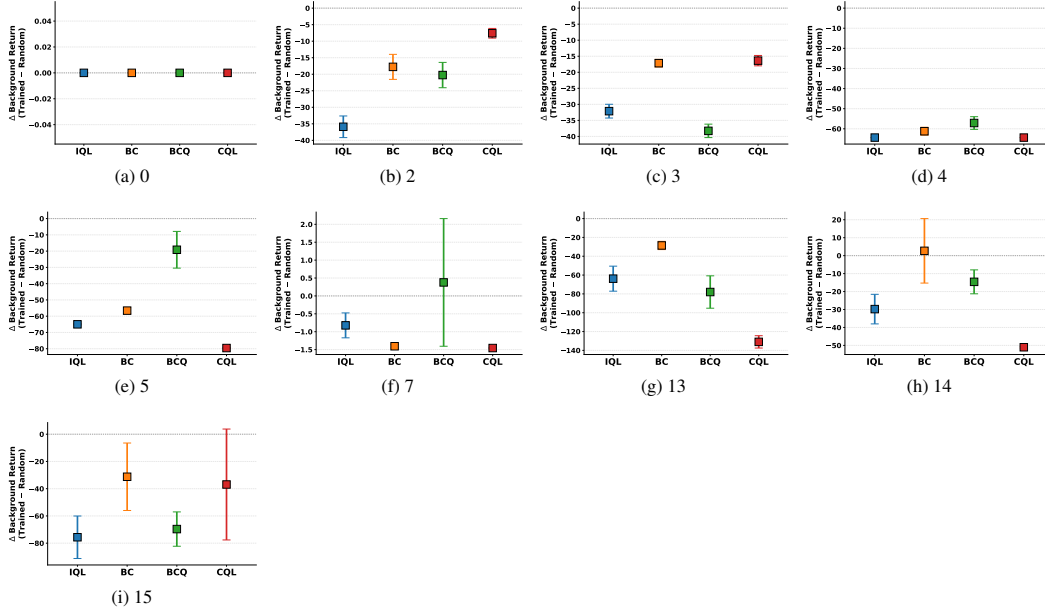


Figure 37: **Clean Up — background-agent return change.** For each scenario in Clean Up, one marker per algorithm showing the change in mean background-agent episode return induced by the trained focal policy, relative to a uniform-random focal policy ( $\Delta = R_{\text{trained}} - R_{\text{random}}$ ). Background returns are read from the in-training eval logs. Negative values indicate the focal policy hurt the background population (exploitation); positive values indicate it helped them.

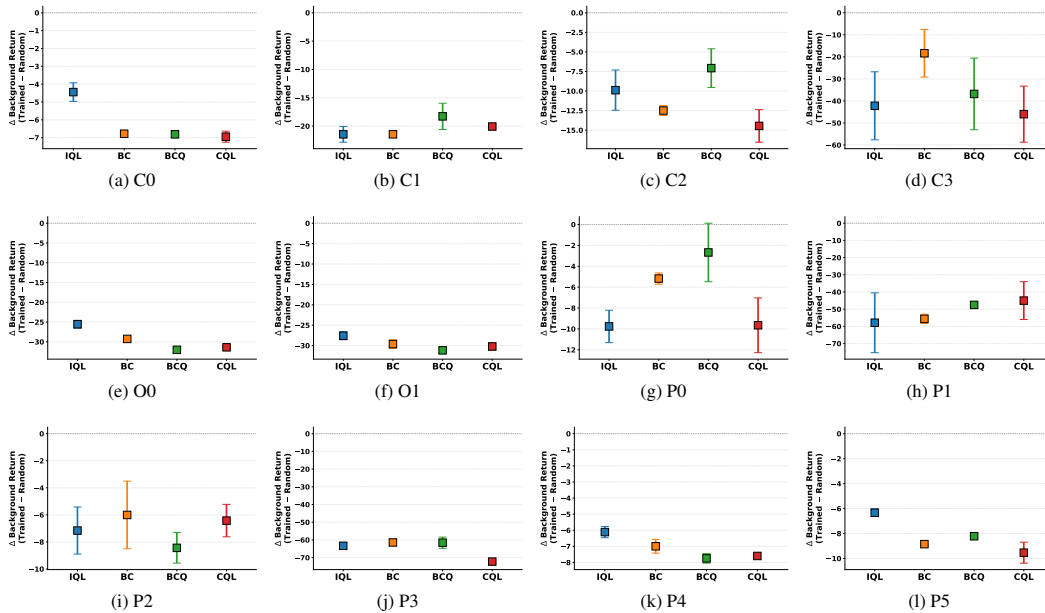


Figure 38: **Commons Harvest — background-agent return change.** For each scenario in Commons Harvest, one marker per algorithm showing the change in mean background-agent episode return induced by the trained focal policy, relative to a uniform-random focal policy ( $\Delta = R_{\text{trained}} - R_{\text{random}}$ ). Background returns are read from the in-training eval logs. Negative values indicate the focal policy hurt the background population (exploitation); positive values indicate it helped them.

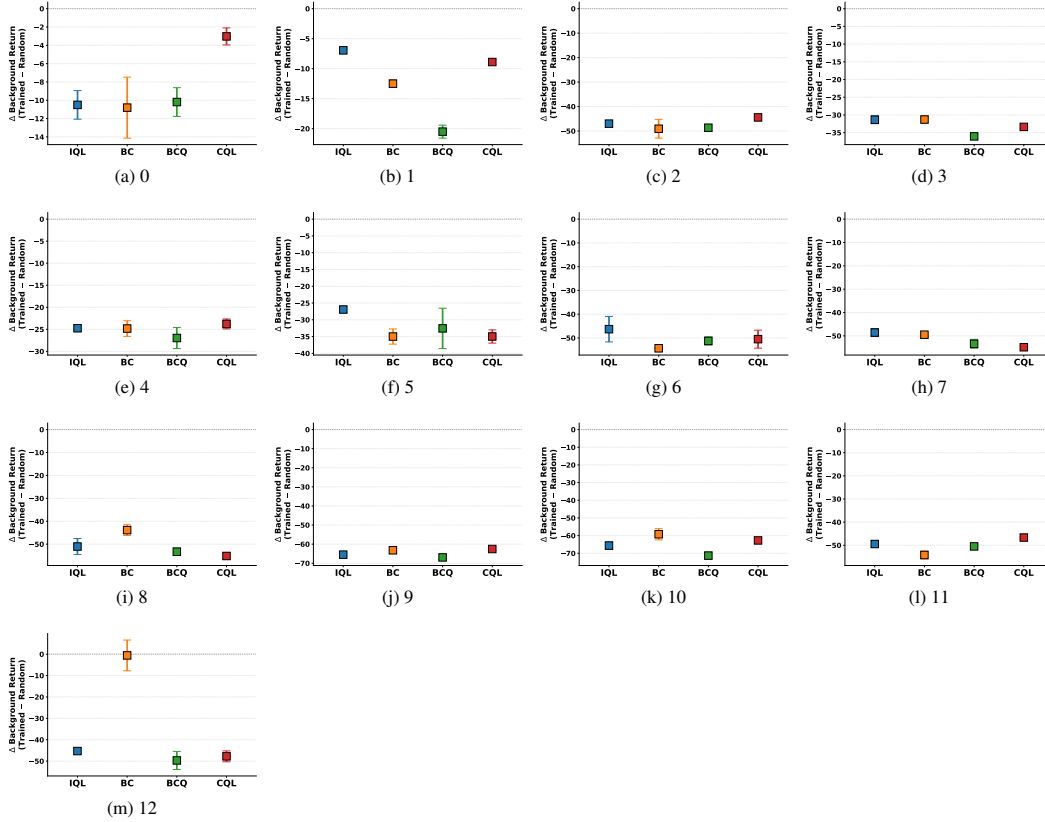


Figure 39: **Allelopathic Harvest — background-agent return change.** For each scenario in Allelopathic Harvest, one marker per algorithm showing the change in mean background-agent episode return induced by the trained focal policy, relative to a uniform-random focal policy ( $\Delta = R_{\text{trained}} - R_{\text{random}}$ ). Background returns are read from the in-training eval logs. Negative values indicate the focal policy hurt the background population (exploitation); positive values indicate it helped them.

525 **G Cross-substrate aggregate scores**

526 Table 7 reports the full set of reliable headline estimators — Median, interquartile mean (IQM), mean,  
 527 and optimality gap — pooled across every (substrate, scenario, seed) sample for Evaluation Settings 1  
 528 and 2. The compact optimality-gap view in the main text (Table 1) is the last row of each block here.

Table 7: **Cross-substrate aggregate scores — Evaluation 1 vs Evaluation 2.** Median, interquartile mean (IQM), mean, and optimality gap of the normalised final-eval return pooled across every (substrate, scenario, seed) sample, with stratified-bootstrap 95% CIs in brackets (reliable). Higher is better for the first three rows; optimality gap ( $\downarrow$ ). Best per row in bold.

<b>Metric</b>	<b>IQL</b>	<b>BC</b>	<b>BCQ</b>	<b>CQL</b>
<i>Evaluation Setting 1 (single-scenario)</i>				
Median	0.75 [0.73, 0.77]	0.72 [0.68, 0.73]	<b>0.80 [0.77, 0.83]</b>	0.63 [0.60, 0.66]
IQM	0.73 [0.71, 0.74]	0.69 [0.68, 0.71]	<b>0.80 [0.78, 0.82]</b>	0.63 [0.62, 0.65]
Mean	0.78 [0.75, 0.81]	0.66 [0.65, 0.67]	<b>0.89 [0.86, 0.92]</b>	0.58 [0.56, 0.59]
Optimality Gap	0.30 [0.28, 0.33]	0.35 [0.33, 0.36]	<b>0.26 [0.24, 0.27]</b>	0.43 [0.42, 0.44]
<i>Evaluation Setting 2 (multi-scenario, scenario label withheld)</i>				
Median	0.72 [0.71, 0.74]	0.68 [0.65, 0.69]	<b>0.73 [0.70, 0.78]</b>	0.58 [0.52, 0.63]
IQM	<b>0.73 [0.72, 0.75]</b>	0.66 [0.65, 0.67]	<b>0.73 [0.71, 0.77]</b>	0.57 [0.54, 0.58]
Mean	0.78 [0.77, 0.80]	0.66 [0.65, 0.68]	<b>0.79 [0.75, 0.83]</b>	0.55 [0.53, 0.56]
Optimality Gap	0.31 [0.29, 0.32]	0.36 [0.35, 0.37]	<b>0.29 [0.27, 0.31]</b>	0.49 [0.48, 0.51]

## 529 **NeurIPS Paper Checklist**

### 530 **1. Claims**

531 Question: Do the main claims made in the abstract and introduction accurately reflect the  
532 paper’s contributions and scope?

533 Answer: [\[Yes\]](#)

534 Justification: The abstract and section 1 describe the protocol, datasets, three evaluation  
535 settings, and benchmarking of four offline RL algorithms; these are delivered in section 4  
536 and section 5 respectively.

### 537 **2. Limitations**

538 Question: Does the paper discuss the limitations of the work performed by the authors?

539 Answer: [\[Yes\]](#)

540 Justification: subsection 4.2 notes that poor Setting 2 performance can stem from offline-RL  
541 pathologies as well as missing partner inference, and section 7 discusses headroom against  
542 the  $p_{90}$  ceiling and the architectural changes likely required to close it. The benchmark is  
543 also restricted to five Melting Pot substrates with discrete actions and pixel observations,  
544 which bounds the generality of the conclusions to real world partner inference.

### 545 **3. Theory assumptions and proofs**

546 Question: For each theoretical result, does the paper provide the full set of assumptions and  
547 a complete (and correct) proof?

548 Answer: [\[Yes\]](#)

549 Justification: The Bayes-optimality argument in section 3 states the latent-context assump-  
550 tion, defines the partner posterior (Eq. 3) and the Bayes-optimal Bellman equation (Eq. 4),  
551 and gives the derivation in full.

### 552 **4. Experimental result reproducibility**

553 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
554 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
555 of the paper (regardless of whether the code and data are provided or not)?

556 Answer: [\[Yes\]](#)

557 Justification: subsection 2.4 describes dataset construction (independent PPO,  $N = 1000$   
558 trajectories of length  $T = 1000$  uniformly across training time), section 5 specifies the  
559 four benchmarked algorithms, recurrent (GRU) architecture, 20,000 gradient updates, and  
560 three seeds per (scenario/substrate/split), and Appendix A lists the substrates, scenarios, and  
561 Setting 3 splits. We will release code.

### 562 **5. Open access to data and code**

563 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
564 tions to faithfully reproduce the main experimental results, as described in supplemental  
565 material?

566 Answer: [\[Yes\]](#)

567 Justification: The Molten Pot datasets (approximately one terabyte of trajectory data across  
568 the 47 scenarios) and benchmarking code will be released alongside the paper under a  
569 permissive open-source license.

### 570 **6. Experimental setting/details**

571 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
572 rameters, how they were chosen, type of optimizer) necessary to understand the results?

573 Answer: [\[Yes\]](#)

574 Justification: section 5 reports the algorithm set, recurrent (GRU) architecture, 20,000  
575 gradient updates, three seeds, and the d4rl-style normalisation; the train/test splits for  
576 Setting 3 are documented in Table 4, and additional hyperparameter detail accompanies the  
577 released code.

### 578 **7. Experiment statistical significance**

579 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
580 information about the statistical significance of the experiments?

581 Answer: [Yes]

582 Justification: Aggregate plots in section 5 use rliable’s stratified-bootstrap 95% confidence  
583 intervals [2, 16], the optimality-gap table (Table 1) reports the same intervals, and per-  
584 scenario figures in the appendix display standard error of the mean across three seeds.

#### 585 8. Experiments compute resources

586 Question: For each experiment, does the paper provide sufficient information on the com-  
587 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
588 the experiments?

589 Answer: [Yes]

590 Justification: Dataset generation and benchmarking compute requirements are documented.

#### 591 9. Code of ethics

592 Question: Does the research conducted in the paper conform, in every respect, with the  
593 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

594 Answer: [Yes]

595 Justification: The work uses simulated multi-agent environments only, involves no human  
596 subjects or personal data, and complies with the NeurIPS Code of Ethics.

#### 597 10. Broader impacts

598 Question: Does the paper discuss both potential positive societal impacts and negative  
599 societal impacts of the work performed?

600 Answer: [N/A]

601 Justification: Molten Pot is an evaluation protocol and dataset for offline social RL on  
602 simulated grid-world substrates; it has no direct deployment pathway and the foreseeable  
603 impact is methodological (improved measurement of social robustness in offline RL).

#### 604 11. Safeguards

605 Question: Does the paper describe safeguards that have been put in place for responsible  
606 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
607 image generators, or scraped datasets)?

608 Answer: [N/A]

609 Justification: The released artefacts are trajectory datasets from simulated grid-world sub-  
610 strates and small offline-RL policy checkpoints; they pose no misuse risk requiring safe-  
611 guards.

#### 612 12. Licenses for existing assets

613 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
614 the paper, properly credited and are the license and terms of use explicitly mentioned and  
615 properly respected?

616 Answer: [Yes]

617 Justification: The substrates build on Melting Pot [1, 25] (Apache 2.0) and we cite the  
618 original substrate sources (Coins [26], Commons Harvest [34]); evaluation utilities use  
619 MARL-eval/rliable [2, 16]; all licenses are respected and credited.

#### 620 13. New assets

621 Question: Are new assets introduced in the paper well documented and is the documentation  
622 provided alongside the assets?

623 Answer: [Yes]

624 Justification: The Molten Pot datasets, splits, and evaluation harness are documented in  
625 section 4 and Appendix A (substrate table, full per-scenario catalogue, and Setting 3 splits),  
626 with per-scenario return histograms in Appendix B; a README accompanies the released  
627 code.

#### 628 14. Crowdsourcing and research with human subjects

629 Question: For crowdsourcing experiments and research with human subjects, does the paper  
630 include the full text of instructions given to participants and screenshots, if applicable, as  
631 well as details about compensation (if any)?

632 Answer: [N/A]

633 Justification: The paper does not involve crowdsourcing or research with human subjects.

634 **15. Institutional review board (IRB) approvals or equivalent for research with human**  
635 **subjects**

636 Question: Does the paper describe potential risks incurred by study participants, whether  
637 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
638 approvals (or an equivalent approval/review based on the requirements of your country or  
639 institution) were obtained?

640 Answer: [N/A]

641 Justification: The paper does not involve research with human subjects, so IRB approval is  
642 not applicable.

643 **16. Declaration of LLM usage**

644 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
645 non-standard component of the core methods in this research? Note that if the LLM is used  
646 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
647 scientific rigor, or originality of the research, declaration is not required.

648 Answer: [N/A]

649 Justification: LLMs are not involved; the benchmarked algorithms (BC, BCQ, IQL, CQL)  
650 and the evaluation protocol do not involve LLMs.