

# Efficient Exploration of Zero-Sum Stochastic Games

Carlos Martin,<sup>1</sup> Tuomas Sandholm<sup>1,2,3,4</sup>

<sup>1</sup>Carnegie Mellon University

<sup>2</sup>Strategy Robot, Inc.

<sup>3</sup>Optimized Markets, Inc.

<sup>4</sup>Strategic Machine, Inc.

cgmartin@cs.cmu.edu, sandholm@cs.cmu.edu

## Abstract

1 We investigate the increasingly important and common game-  
2 solving setting where we do not have an explicit game de-  
3 scription but only oracle access through playing, such as in  
4 financial or military simulations and computer games. Dur-  
5 ing a limited-duration learning phase, the algorithm can con-  
6 trol the actions of both players to try to learn the game and  
7 how to play it well. After that, the algorithm has to pro-  
8 duce a policy that plays well against *any* opponent. For the  
9 stochastic game setting (and normal-form games as a spe-  
10 cial case), we propose using the distribution of state-action  
11 value functions induced by a belief distribution over possi-  
12 ble environments in various ways. We compare the perfor-  
13 mance of many exploration policies for this task, including  
14 generalizations of exploration policies from the single-agent  
15 setting to games, as well as new ones. The multiagent set-  
16 ting raises new challenges and qualitative changes in behav-  
17 ior. Experiments across different environments show orders  
18 of magnitude sample-efficiency improvements over a well-  
19 known model-free approach. Modified versions of Thompson  
20 sampling and UCB are consistently among the most sample  
21 efficient. The former is also computationally efficient.

## Introduction

22 We study how to efficiently explore zero-sum games whose  
23 payoffs and dynamics are initially unknown. The agent is  
24 given a certain number of episodes to learn as much useful  
25 information about the game as possible. During this learn-  
26 ing, the agent can control the play of both players in the  
27 game and the rewards obtained are fictional and thus do not  
28 count toward the evaluation of the final strategy. After this  
29 exploration phase, the agent must recommend a strategy for  
30 playing the game that should be minimally exploitable by  
31 an adversary (who has complete knowledge of the environ-  
32 ment and can thus play optimally against it). This setup is  
33 called *pure exploration* in single-agent *reinforcement learn-*  
34 *ing (RL)*. This setup is important for simulation-based games  
35 in which a black-box simulator is queried with strategies to  
36 obtain samples of the players' resulting utilities (Vorobey-  
37 chik and Wellman 2009), as opposed to the rules of the game  
38 being explicitly given. For example, in many military set-  
39 tings, war game simulators are used to generate strategies,  
40

and then the strategies need to be ready to deploy in case  
of actual war (Marchesi, Trovò, and Gatti 2019). Another  
example is finance, where trading strategies are generated  
in simulation, and then they need to be ready for live trad-  
ing. A third example is video games such as Dota 2 (Berner  
et al. 2019) and Starcraft II (Vinyals et al. 2019), where AIs  
can be trained largely through self play, but most prior tech-  
niques have not focused on exploration approaches that yield  
strategies that have low exploitability against *every possible*  
adversary strategy.

This raises the challenge not only of learning approximate  
equilibria with noisy observations, but also of learning in as  
few queries as possible, since running the simulator is usu-  
ally expensive. Prior work on games such as Dota 2 and Star-  
craft II has not evaluated the exploitability of the learned  
strategies. In this paper, we experiment on games that are  
small enough that game-theoretic exploitability of the rec-  
ommended strategy can be evaluated. We study exploration  
in two-player zero-sum stochastic games, which subsume  
two-player zero-sum extensive-form games that are other-  
wise of perfect-information but have simultaneous moves.  
Our approach is *model driven*: our algorithms incrementally  
build a model of the game and uses it to guide exploration.

## Related research

64 In single-agent RL, Q-learning (Watkins and Dayan 1992)  
65 learns state-action values directly in a model-free way, that  
66 is, without learning the structure of the environment. Dear-  
67 den, Friedman, and Russell (1998) extended it to incor-  
68 porate uncertainty by propagating probability distributions  
69 over the Q values in order to compute a myopic approxi-  
70 mation of the value of information, which measures the ex-  
71 pected improvement in future decision quality that can be  
72 gained from exploration. Bellemare, Dabney, and Munos  
73 (2017) and O'Donoghue et al. (2018) argue for the impor-  
74 tance of the Q value distribution and propose a new version  
75 of Bellman updating that incorporates uncertainty.  
76

A similar problem has been studied in single-agent  
model-based RL. *Posterior sampling RL* (Osband and Roy  
2017; Agrawal and Jia 2017) samples an environment from  
the agent's belief distribution, follows a policy that is op-  
timal with respect to it, and then updates the agent's be-  
liefs about the environment with the resulting observations.  
Zhou, Li, and Zhu (2020) extended it to zero-sum imperfect-

84 information extensive-form games and presented an algo-  
85 rithm that converges to a Nash equilibrium at a bounded rate.

86 Even in single-agent settings, sampling a new policy on  
87 every step within an episode is inefficient because it does  
88 not do *deep exploration*, which accounts not only for in-  
89 formation gained by taking an action but also for how the  
90 action may position the agent to more effectively acquire in-  
91 formation later (Russo et al. 2018). In deep exploration, a  
92 single policy is chosen at the beginning of each episode and  
93 followed for its duration. One approach to deep exploration  
94 chooses actions that are optimal with respect to a value func-  
95 tion that is sampled from an ensemble (Osband et al. 2016,  
96 2019). Each element of the ensemble is a deep neural net-  
97 work trained with deep Q-learning (Mnih et al. 2015), and  
98 the ensemble constitutes a belief distribution over possible  
99 value functions of the environment. It incentivizes exper-  
100 imentation with actions of uncertain value because uncer-  
101 tainty induces variance in the sampled value estimate. Chen  
102 et al. (2017) also use an ensemble of Q functions but, instead  
103 of sampling from them, use the resulting upper confidence  
104 bounds. Mavrin et al. (2019) combine a decaying sched-  
105 ule with exploration bonuses computed from upper quan-  
106 tiles of the learned distribution. Littman (1994) describes  
107 a Q-learning-like algorithm for finding optimal policies for  
108 one player in stochastic games (Shapley 1953) when playing  
109 against an opponent that the algorithm does not control.

110 Sandholm and Crites (1996) study RL in a repeated game.  
111 They study the role of other agents making the setting  
112 stochastic for a learner, the role of exploration, and con-  
113 vergence to cycles of different lengths, and how recurrent  
114 neural networks can, in principle, help with those issues.  
115 Claus and Boutilier (1998) study RL in cooperative set-  
116 tings, showing that several optimistic exploration strategies  
117 increase the likelihood of reaching an optimal equilibrium.  
118 Wang and Sandholm (2002) describe an algorithm that con-  
119 verges almost surely to an optimal equilibrium in any team  
120 stochastic game. Conitzer and Sandholm (2003) present BL-  
121 WoLF, a framework for learnability in repeated zero-sum  
122 games where the cost of learning is measured by the losses  
123 the learning agent accrues, and guaranteed learnability re-  
124 sults for families of games. Hu and Wellman (2003) present  
125 Nash Q-learning for general-sum stochastic games, which is  
126 guaranteed to converge to an equilibrium if all agents fol-  
127 low the algorithm and the stage games satisfy certain highly  
128 restrictive conditions. Ganzfried and Sandholm (2009) de-  
129 sign algorithms for computing equilibria in special classes of  
130 stochastic games of imperfect information. Heinrich and Sil-  
131 ver (2016) introduced an approach to learning approximate  
132 Nash equilibria without prior domain knowledge by combin-  
133 ing fictitious self-play with deep RL. Lanctot et al. (2017)  
134 introduce a meta-algorithm for *multiagent RL (MARL)* based  
135 on approximate best responses to mixtures of policies gener-  
136 ated using deep RL. Srinivasan et al. (2018) apply candi-  
137 date policy update rules to model-free MARL in adver-  
138 sarial sequential decision problems, showing empirical con-  
139 vergence to approximate Nash equilibria in self play. Cas-  
140 grain, Ning, and Jaimungal (2019) study Nash Q-learning  
141 but they use a neural network to model the Q function, de-  
142 compositing it into a sum of the state value function and a

specific form of action advantage function. Sokota, Ho, and  
Wiedenbeck (2019) use neural networks to learn a mapping  
from mixed-strategy profiles to deviation payoffs in order to  
approximate role-symmetric equilibria in large simulation-  
based games. Lockhart et al. (2019) present exploitability  
descent, an algorithm to compute approximate equilibria in  
extensive-form games by direct policy optimization against  
worst-case opponents.

Our work also relates to the multi-armed bandit problem.  
In the standard version, the agent must—over a given period  
of time—acquire information about the mean payoff of each  
action while simultaneously trying to maximize the cumu-  
lative payoff. Because of this tradeoff, multi-armed bandit  
problems exemplify the exploration-exploitation dilemma.  
In the *pure exploration* version, there is a learning phase  
first, during which the rewards obtained are fictional and do  
not count toward the evaluation (Bubeck, Munos, and Stoltz  
2009). Then the agent recommends an arm (that is, action)  
to play going forward. The agent’s performance is measured  
purely by the effectiveness of her recommended action. This  
performance measure is called *simple regret*, in contrast to  
the cumulative regret of the standard problem where rewards  
throughout the process count. Our paper focuses on the mul-  
tiagent generalization of pure exploration.

Garivier, Kaufmann, and Koolen (2016) study the prob-  
lem of pure exploration in the context of a *sequential-move*  
game with the aim of identifying an  $\varepsilon$ -maxmin action with  
probability at least  $1 - \delta$ . Marchesi, Trovò, and Gatti (2019)  
tackle the problem of learning equilibria in simulation-based  
games of infinite strategy spaces with high confidence using  
as few simulator queries as possible. They propose an algo-  
rithm for the fixed-confidence setting (guaranteeing the  
desired confidence level while minimizing the number of  
queries) and one for the fixed-budget setting (maximizing  
the confidence without exceeding the given maximum num-  
ber of queries).

## Zero-sum stochastic games

$\Delta\mathcal{X} \subseteq \mathcal{X} \rightarrow \mathbb{R}$  denotes the set of all probability distribu-  
tions on a set  $\mathcal{X}$ .  $\mathcal{N}$  denotes a set of players.  $i \in \mathcal{N}$  denotes a  
player.  $-i$  denotes the remaining player(s).  $\mathcal{S}$  denotes a set  
of states.  $s \in \mathcal{S}$  denotes a state.  $\mathcal{A}_i(s)$  denotes  $i$ ’s set of ac-  
tions in  $s$ .  $a_i \in \mathcal{A}_i(s)$  denotes  $i$ ’s action in  $s$ .  $a \in \prod_i \mathcal{A}_i(s)$   
denotes an action profile in  $s$ .  $\sigma_i \in \Delta\mathcal{A}_i(s)$  denotes  $i$ ’s strat-  
egy in  $s$ .  $\sigma \in \prod_i \Delta\mathcal{A}_i(s)$  denotes a strategy profile in  $s$ .  
 $\pi_i \in \prod_s \Delta\mathcal{A}_i(s)$  denotes  $i$ ’s policy.  $\pi \in \prod_i \prod_s \Delta\mathcal{A}_i(s)$  de-  
notes a policy profile.  $r_i \in \prod_s (\prod_i \mathcal{A}_i(s) \rightarrow \mathbb{R})$  denotes  $i$ ’s  
reward function.  $\delta \in \prod_s (\prod_i \mathcal{A}_i(s) \rightarrow \Delta\mathcal{S})$  denotes a tran-  
sition function.  $\gamma \in [0, 1]$  denotes a discount factor.  $s_{\text{init}} \in \mathcal{S}$   
denotes the initial state.

A stochastic game is a tuple  $(\mathcal{N}, \mathcal{S}, \mathcal{A}, r, \delta, \gamma, s_{\text{init}})$ . We  
consider *finite-horizon* games, so  $\gamma = 1$ , the transition graph  
induced by  $\delta$  is acyclic, and terminal states have  $\mathcal{A}_i(s) = \emptyset$ .  
The game proceeds as a sequence of *stage games* in which  
all players choose actions simultaneously, each player re-  
ceives a reward, and the game transitions to a new state. This  
process repeats until a terminal state is reached. We consider  
two-player zero-sum games, so  $\mathcal{N} = \{1, 2\}$  and  $\sum_i r_i = 0$ .

This model is quite general and captures many types of

201 games. If there is a single nonterminal state, we have a matrix  
 202 game (Sandholm and Crites 1996). If  $|\mathcal{A}_{-i}(s)| = 1$  for  
 203 every  $s \in \mathcal{S}$ , we have a Markov decision process (MDP) for  $i$ .  
 204 If both of the above conditions hold, we have a multi-armed  
 205 bandit. If the state transition graph induced by  $\delta$  is a tree, we  
 206 have a perfect-information extensive-form game.

## 207 Policy recommendation under uncertainty

Let  $u_i(g, \pi)$  denote the expected return to player  $i$  of policy profile  $\pi$  under game  $g$ . Suppose we face the task of recommending a policy for player  $i$ , under the assumption that our opponent will play optimally against it. To maximize our expected utility, we should recommend

$$\operatorname{argmax}_{\pi_i} \min_{\pi_{-i}} u_i(g, \pi) \quad (1)$$

We define the *regret* incurred by  $\hat{\pi}_i$  under  $g$  as

$$\max_{\pi_i} \min_{\pi_{-i}} u(g, \pi) - \min_{\pi_{-i}} u(g, \hat{\pi}_i, \pi_{-i}) \quad (2)$$

208 It measures how exploitable the recommended policy  $\hat{\pi}_i$  is  
 209 in comparison to the optimal policy for  $g$ .

But suppose that we are uncertain about  $g$ . More precisely, suppose the reward and transition functions of  $g$  are parameterized by some unknown parameter  $\theta \in \Theta$ , where our beliefs about  $\theta$  are modelled by some distribution  $D_\Theta : \Delta\Theta$  that induces a corresponding distribution  $D : \Delta\mathcal{G}$  over games (we call this the *belief distribution*). In that case, to maximize our expected utility, we should recommend

$$\operatorname{argmax}_{\pi_i} \mathbb{E}_{g \sim D} \min_{\pi_{-i}} u_i(g, \pi) \quad (3)$$

Unfortunately, this typically has no closed-form solution. Instead, we replace the mean with a sample mean:

$$\hat{\pi}_i = \operatorname{argmax}_{\pi_i} \sum_j \min_{\pi_{-i}} u_i(g_j, \pi) \quad (4)$$

where  $g_j \sim D$  for  $j \in \mathcal{J}$ . That is, we find a policy that optimizes over an *ensemble* of games, rather than a single game. The probability that the optimality gap

$$\max_{\pi_i} \mathbb{E}_{g \sim D} \min_{\pi_{-i}} u_i(g, \pi) - \mathbb{E}_{g \sim D} \min_{\pi_{-i}} u_i(g, \hat{\pi}_i, \pi_{-i}) \quad (5)$$

210 exceeds  $\varepsilon \geq 0$  is at most  $e^{-2|\mathcal{J}|\varepsilon^2/\Delta^2}$ , where  $\Delta$  is  $i$ 's payoff  
 211 range (maximum possible payoff minus minimum possible  
 212 payoff). See the appendix for a proof.

## 213 Computation

We can compute this policy through backward induction. Let

$$v : \prod_j \prod_s \mathbb{R} \quad q : \prod_j \prod_s \prod_a \mathbb{R} \quad \pi_i : \prod_s \Delta\mathcal{A}_i(s) \quad (6)$$

where  $v$  yields a value for each state in each game,  $q$  yields a value for each action profile in each state in each game, and  $\pi_i$  is a policy for  $i$ . Let  $r_j$  and  $\delta_j$  be the reward and transition functions of  $g_j$ . Then let

$$v_j(s) = \min_{a_{-i}} q_j(s, \pi_i(s), a_{-i}) \quad (7)$$

$$q_j(s, \sigma_i, a_{-i}) = \sum_{a_i} \sigma_i(a_i) q_j(s, a) \quad (8)$$

$$q_j(s, a) = r_{j,i}(s, a) + \sum_{s'} \delta_j(s, a, s') v_j(s') \quad (9)$$

$$\pi_i(s) = \operatorname{argmax}_{\sigma_i} \sum_j \min_{a_{-i}} q_j(s, \sigma_i, a_{-i}) \quad (10)$$

where  $v_j(s) \stackrel{\text{def}}{=} 0$  for terminal states. To compute the first and last lines, we formulate and solve the linear program

$$\text{maximize } \mathbf{1} \cdot \alpha \quad \text{over } \alpha : \mathbb{R}^{\mathcal{J}} \quad \beta : \mathbb{R}^{\mathcal{A}_i(s)} \quad (11)$$

$$\text{subject to } \mathbf{1} \cdot \beta = 1 \quad (12)$$

$$\beta \geq \mathbf{0} \quad (13)$$

$$\alpha \leq \Gamma(a_{-i}) \cdot \beta \quad \forall a_{-i} : \mathcal{A}_{-i}(s) \quad (14)$$

214 where  $\Gamma(a_{-i})_j = q_j(s, a_{-i})$ , and let  $v_j(s) = \alpha_j, \pi_i(s) = \beta$ .  
 215 The dual variables associated with the inequality constraints  
 216 contain best responses from player  $-i$ , that is,  $\pi_{-i,j}(s)$ .

217 The recommended policy is a non-local function of both  
 218 observations and the initial belief distribution. Depending on  
 219 the belief distribution, there may be nontrivial correlations  
 220 across states, including unvisited states. For example, if we  
 221 believe there is a treasure chest at location A or B (but not  
 222 both), then once we have explored A and found nothing, we  
 223 will know that the treasure chest is at B and behave accordingly,  
 224 despite not having visited B. This example illustrates the  
 225 advantage of our Bayesian approach, which allows one  
 226 to incorporate any useful domain-specific knowledge into  
 227 the initial belief distribution.

## 228 Exploration policies for games

229 Our problem is as follows: We face a stochastic game with  
 230 unknown rewards and/or transitions. Our beliefs about the  
 231 true game are modelled by a distribution  $D$  over possible  
 232 games, from which we have the ability to sample. For each  
 233 of  $T$  episodes, we choose a *policy profile*, observe an episode  
 234 of gameplay under that policy profile—including the transi-  
 235 tions and (possibly stochastic) rewards on each timestep—  
 236 and update our beliefs according to these observations. After  
 237  $T$  episodes, we recommend a policy for Player  $i$  (see Equa-  
 238 tion 4). This leaves the question of how to explore over the  
 239 course of the  $T$  episodes, since this determines how useful  
 240 the beliefs we end up with will be. In this section, we de-  
 241 sign a number of exploration policies, and later in the paper  
 242 compare them experimentally.

Let  $v_t(D)$  denote the expected utility of a belief distribution  $D$  when  $t$  exploration episodes remain. Then

$$v : \mathbb{N} \rightarrow \Delta\mathcal{G} \rightarrow \mathbb{R} \quad (15)$$

$$v_0(D) = \max_{\pi_i} \mathbb{E}_{g \sim D} \min_{\pi_{-i}} u_i(g, \pi) \quad (16)$$

$$v_{t+1}(D) = \max_{\pi} \mathbb{E}_{D'|D, \pi} v_t(D') \quad (17)$$

243 where  $D'|D, \pi$  denotes the belief distribution we acquire af-  
 244 ter an episode, conditioned on the fact that  $D$  is our belief  
 245 distribution before the episode and  $\pi$  is the policy profile we  
 246 use during the episode.

247 Thus, if our current belief distribution is  $D$  and there are  $t$   
 248 remaining exploration episodes, the optimal policy profile to  
 249 use for exploration is  $\operatorname{argmax}_{\pi} \mathbb{E}_{D'|D, \pi} v_t(D')$ , which can  
 250 be computed recursively as described above. However, this  
 251 computation is intractable. The nested maximizations form  
 252 a tree whose depth is the number of remaining episodes, and  
 253 the branching factor (the number of policy profiles to opti-  
 254 mize over in each maximization) is itself exponential in the  
 255 size of the game. Even computing the optimal exploration  
 256 policy with a single-step look-ahead (that is, letting  $t = 1$ )  
 257 is intractable, for the latter reason.

258 In the rest of the paper, we will design exploration policies  
 259 that are more computationally feasible. Our first policy is the  
 260 **uniform** policy, which simply selects actions for all players  
 261 uniformly at random:  $\pi_i(s) = \operatorname{uniform}(\mathcal{A}_i(s))$ . The **min-**  
 262 **count** policy selects in each state the action profile that has  
 263 been chosen the least so far:  $a = \operatorname{argmin}_a \#(s, a)$ , where  
 264  $\#(s, a)$  is the number of times  $a$  has been chosen in  $s$ . Ties  
 265 are broken uniformly at random. The **greedy** policy uses the  
 266  $\operatorname{argmax}$ -sum-min policy given by Equation 4 both players.

267 The main idea behind this policy (as well as other policies  
 268 described below) is that, if we are reasonably certain about  
 269 where the Nash equilibrium of the game is, it makes sense to  
 270 try to refine our knowledge by staying near this equilibrium,  
 271 rather than waste time on game trajectories that we believe  
 272 are far away from equilibrium gameplay.

273 The greedy policy, at least in the limit of infinite samples  
 274 from the belief distribution, is not guaranteed to converge to  
 275 correct beliefs about the game, because it may get stuck with  
 276 false max-mean-min policies for both players from which it  
 277 is continuous to sample forever, ignoring the rest of the game.

278 To solve this problem, we can force the greedy policy to  
 279 explore more by combining it with the uniform policy, yield-  
 280 ing the  $\varepsilon$ -**greedy** policy. More precisely, at each individual  
 281 timestep,  $\varepsilon$ -greedy chooses to execute the uniform policy in-  
 282 stead of the greedy policy with probability  $\varepsilon$ .

283 Our next exploration policy is inspired by Thompson sam-  
 284 pling (Thompson 1933), a remarkably simple successful  
 285 heuristic for solving exploration-exploitation dilemmas in  
 286 multi-armed bandit problems (Ortega and Braun 2010). It  
 287 consists of playing an action according to the probability that  
 288 it is the optimal action. This is done by sampling a belief  
 289 from the belief distribution and then acting optimally with  
 290 respect to that belief. It has become a popular approach for  
 291 RL (Russo et al. 2018), and convergence results have been  
 292 obtained that show it is asymptotically optimal and well be-  
 293 haved (Osband and Roy 2017; Agrawal and Jia 2017).

294 Our **Thompson policy**, a generalization of Thompson  
 295 sampling to games, samples a single game from our belief  
 296 distribution at the beginning of each episode, and then uses  
 297 the  $\operatorname{argmax}$ -min policies (Equation 1) for both players in that  
 298 game for the duration of the episode. Like the greedy explo-  
 299 ration policy, it biases gameplay toward what is believed to  
 300 be the Nash equilibrium, but unlike greedy exploration, it  
 301 leaves more room for exploration according to how uncer-  
 302 tain the agent is. The more certain we are of the true game,  
 303 and thus of where its Nash equilibrium lies, the more often  
 304 we will sample trajectories played by that equilibrium, re-  
 305 fining our knowledge of that region of the game.

306 Recall that our ultimate goal is to recommend a policy for  
 307 player  $i$ , assuming  $-i$  will play optimally against it. Because  
 308 of this, we also propose the **semi-Thompson** policy, which  
 309 uses  $-i$ 's *best response* to  $i$ 's  $\operatorname{argmax}$ -min policy, rather than  
 310  $-i$ 's own  $\operatorname{argmax}$ -min policy.

Auer, Cesa-Bianchi, and Fischer (2002) introduced the  
 UCB1 algorithm for multi-armed bandit problems, proving  
 that it achieves optimal regret up to a multiplicative constant.  
 UCB1 selects the action  $a$  with the highest *upper confidence*  
*bound* (UCB) on its utility, where (for utilities in  $[0, 1]$ ),

$$\operatorname{UCB1} f(g, a) = \mathbb{E}_{g \sim D} f(g, a) + \sqrt{\frac{2 \ln \sum_{a'} \#(a')}{\#(a)}} \quad (18)$$

This type of approach is known as ‘‘optimism in the face of  
 uncertainty’’, and it encourages exploration of actions with  
 more uncertain payoffs. In theory, the analogous policy for  
 our setting would be

$$\operatorname{argmax}_{\pi_i} \operatorname{UCB1} \min_{g \sim D} \min_{\pi_{-i}} u_i(g, \pi) \quad (19)$$

311 Unfortunately, we do not have a tractable way to compute  
 312 this  $\pi_i$ . Instead, we use the following approximation:

313 First, we compute the value  $Q_i(g, s, a)$  of each state and  
 314 action profile for each *individual* game  $g$  under Nash equilib-  
 315 rium gameplay *in that game*. To do this, we reuse the back-  
 316 ward induction algorithm but let  $\mathcal{J} = \{j\}$  be a singleton,  
 317 where  $g_j$  is the individual game under consideration.

Having done this, our **UCB1** policy is

$$\pi_i(s) = \operatorname{argmax}_{\sigma_i} \operatorname{UCB1} \min_{g \sim D} \min_{\sigma_{-i}} Q_i(g, s, \sigma), \text{ where} \quad (20)$$

$$\operatorname{UCB1} f(g, \sigma_i) = \mathbb{E}_{g \sim D} f(g, \sigma_i) + \sqrt{\frac{2 \ln \sum_{a_i} \#(a_i)}{\#(\sigma_i(a_i))}} \quad (21)$$

where we define an analog of action counts for mixed strate-  
 gies, consisting of their expected action count:

$$\#(\sigma_i(a_i)) = \sum_{a_i} \sigma_i(a_i) \#(a_i) \quad (22)$$

We approximate the mean with an *empirical* mean over the  
 samples from our belief distribution,  $\{g_j\}_{j \in \mathcal{J}}$ .

$$\mathbb{E}_{g \sim D} f(g, \sigma_i) \approx \frac{1}{|\mathcal{J}|} \sum_j f(g_j, \sigma_i) \quad (23)$$

Kaufmann, Cappe, and Garivier (2012) introduced a  
 Bayesian version of UCB and proved that it satisfies finite-  
 time regret bounds that imply asymptotic optimality. It se-  
 lects the action whose mean payoff  $1 - \frac{1}{t}$  quantile is highest,  
 where  $t$  is the current exploration episode. Our **Bayes-UCB**  
 policy is the same as the UCB1 policy, except that we mod-  
 ify the UCB:

$$\operatorname{UCB1} f(g, \sigma_i) = \operatorname{quantile}_{g \sim D} f(g, \sigma_i) \quad (24)$$

318 We approximate the quantile with a linearly-interpolated  
 319 *empirical* quantile over samples from our belief distribution.

We then observed that, in a zero-sum game, it is not possible to be optimistic about payoffs for both players simultaneously in an internally self-consistent way—because if the payoff is high for one, it must be low for the other. Thus, for zero-sum games, UCB-like policies like the above are conceptually questionable, because the optimistic assumptions they make seem to be inconsistent. Therefore, we introduce an exploration policy that we call *exploitability optimism*. It is motivated by the success of optimism in exploration in single-agent settings, but does not use the inconsistent optimism about the players’ payoffs. Rather, it samples an action profile from the strategy profile that has the lowest lower confidence bound (LCB) on the players’ combined *exploitabilities*. More precisely, our **sum-exploit** and **max-exploit** policies are

$$\pi(s) = \operatorname{argmin}_{\sigma} \operatorname{LCB}_{g \sim D} f(\operatorname{expl}(Q(g, s), \sigma)) \quad (25)$$

for  $f(x) = \sum_i x_i$  and  $f(x) = \max_i x_i$  respectively, where

$$\operatorname{expl}(u, \sigma)_i = \max_{\sigma'_i} u_i(\sigma'_i, \sigma_{-i}) - u_i(\sigma) \quad (26)$$

learning phase. We experiment in small-to-medium-sized games in order to be able to evaluate the regret exactly.

We use two kinds of benchmarks. The first kind consists of randomly-generated normal-form games with binary  $\{0, 1\}$  (that is, Bernoulli-distributed) rewards. We assume the agent does not know the mean reward (that is, the Bernoulli parameter) associated with any action profile. The agent starts with uniform priors (which are a special case of beta priors) over the parameter for each action profile. Because the beta distribution is a conjugate prior for the Bernoulli distribution, we can easily update these priors according to incoming observations—by simply incrementing the pseudocount corresponding to the observed reward.

Figures 1 and 4 illustrates the performance of different exploration policies on such  $10 \times 10$  normal-form games. We used 100 trials. For the policies requiring it, 100 game samples and 100 strategy or strategy profile samples were used. For the  $\varepsilon$ -greedy exploration policy we used  $\varepsilon = 0.1$ . Our Python 3.8 implementation used Gurobi 9.1 and ran on a 2.5 GHz Quad-Core Intel Core i7 processor.

measures the suboptimality of player  $i$ ’s strategy under payoffs  $u$  and strategy profile  $\sigma$ . As before, we approximate the LCB with a linearly-interpolated empirical quantile—this time of  $\frac{1}{t}$ , where  $t$  is the current exploration episode—over the samples from our belief distribution,  $\{g_j\}_{j:\mathcal{J}}$ .

**Strategy optimization.** For the exploration policies that require optimization over a space of strategies or strategy profiles and lack a closed-form solution—namely UCB1, Bayes-UCB, sum-exploit, and max-exploit—we replace optimization over the uncountable space  $\Delta \mathcal{A}_i(s)$  with optimization over a finite subset  $\mathcal{X}_i \subseteq \Delta \mathcal{A}_i(s)$  of strategies sampled uniformly at random from  $\Delta \mathcal{A}_i(s)$ .

**Convergence.** In our experiments, presented later, all these exploration policies appeared to converge to zero regret. Even the greedy exploration policy did not get stuck because we use a stochastic estimator based on samples from  $D$ . Convergence is *guaranteed* if throughout the exploration process, for each action profile, there is a nonzero probability of sampling a game  $g$  such that the algorithm-selected Nash equilibrium of  $g$  plays that action profile with nonzero probability. This, in turn, can be *guaranteed* by adding a small amount of uniform randomness to any of the exploration policies, as we did with  $\varepsilon$ -greedy. (In our experiments, we did not do that, because they converged anyway.)

The key question is how quickly these exploration policies converge to low-exploitability strategies. Although establishing theoretical worst-case bounds of convergence for our methods would be interesting future work, methods with the best theoretical bounds are not necessarily those that perform best in practice—as the single-agent reinforcement learning literature shows. For that reason, we consider it important to investigate how quickly our methods converge *in practice*.

## Experiments

We compared the performance of the exploration policies we proposed. We measure how quickly the regret of the recommended policy decreases as a function of episodes in the

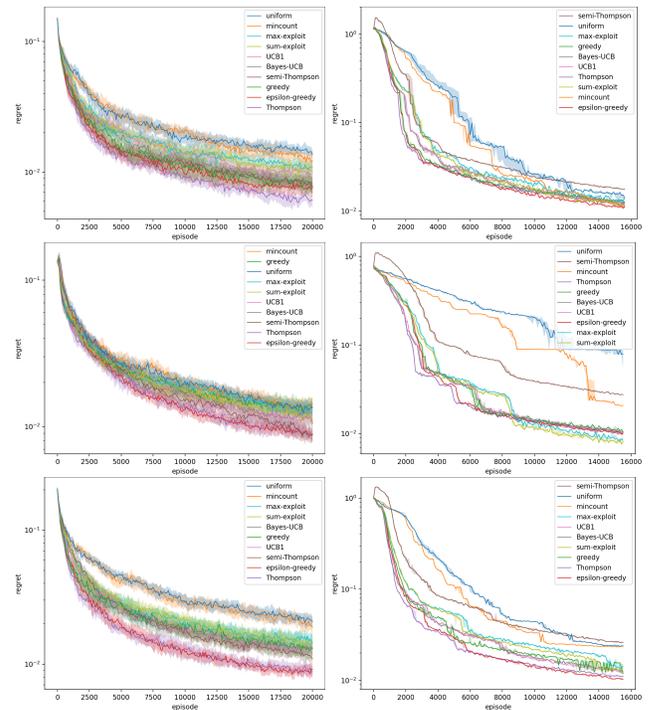


Figure 1: Exploration policies on normal-form games (left) and bombardment games (right). See appendix for more figures.

As expected, the naive exploration policies—uniform and mincount—underperform the more sophisticated exploration policies. The  $\varepsilon$ -greedy policy attains better performance than the greedy policy. The UCB1, Bayes-UCB, and exploitability policies also perform well. Our Thompson and semi-Thompson policies are consistently among the best-performing policies. Thompson has an advantage in terms of computational efficiency (see Figure 2) because it only needs to sample and solve *one* game from its belief distribu-

386 tion.

387 One can generalize this to multi-state stochastic games  
388 with unknown transition probabilities by using a Dirichlet  
389 distribution (which is the conjugate prior of the categorical  
390 distribution) for each state and action profile, incrementing  
391 the pseudocount corresponding to the observed next state.  
392 Our second kind of benchmark studies stochastic games  
393 where the rewards are known but the transition probabili-  
394 ties are not. To study this setting, we created a benchmark  
395 called the *bombardment game*, illustrated in Figure 3. It  
396 is a gridworld-like setting which starts with Player 1 (the  
397 solid disk) in the top-left corner and ends when they reach  
398 the bottom-right corner. In each time step, Player 1 stays  
399 put or chooses one of the 4 cardinal directions to move in.  
400 Player 2 (the ring) simultaneously chooses to bomb either  
401 Player 1’s current position or one of its 4 neighboring posi-  
402 tions. Player 1’s objective is to minimize their damage. They  
403 accomplish this by reaching the goal as quickly as possible  
404 while remaining unpredictable enough to dodge Player 2’s  
405 bombardments to an optimal extent. Player 1 receives a re-  
406 ward of  $-1$  whenever their next position coincides with  
407 Player 2’s ring. The layouts we tested are shown in Figure 3.  
408 Each episode had a 40-step horizon.

409 Player 1 learned to prefer areas with more room for ma-  
410 neuvering and proximity to the goal. Player 2 knows this and  
411 learned its bombardment policy accordingly. This interplay  
412 resulted in complex emergent behavior. The performance of  
413 the different exploration policies on the gap, nest, and dou-  
414 ble fork layouts is shown in Figure 5.

415 Compared to the normal-form game setting, we observe  
416 some significant differences in the performance of the ex-  
417 ploration policies. One important difference is that the semi-  
418 Thompson policy consistently performs worse than almost  
419 all the other policies, and asymptotically seems to be out-  
420 paced by them. This might be due to the fact that it does  
421 not perform enough exploration because it uses only the de-  
422 terministic best response for Player 2 to Player 1’s policy  
423 (under the sampled game). Curiously, its regret increases  
424 sharply in the initial phase of exploration. Another impor-  
425 tant difference is that there is a much larger gap between  
426 the naive exploration policies (uniform and mincount) and  
427 the other exploration policies than in the normal-form game  
428 setting, at least initially. This suggests that, in larger and  
429 more complex environments, there is a bigger advantage  
430 in adopting more sophisticated exploration policies. Finally,  
431 our Thompson policy is still consistently among the best-  
432 performing policies, but the exploitability policies tend to  
433 outperform the rest more frequently than they do in the  
434 normal-form game setting. For all of the exploration poli-  
435 cies except mincount, uniform, greedy, and  $\epsilon$ -greedy, the dy-  
436 namic program across the states is the same. The computa-  
437 tion time scales linearly with the number of states and with  
438 the length of the horizon. The relative run time of the dif-  
439 ferent policies is thus the same as in the normal-form case  
440 shown in Figure 2. Therefore, our Thompson policy can be  
441 recommended due to both its solution quality and its speed  
442 advantage.

443 It is not necessary to start with a uniform prior. For ex-  
444 ample, we may be certain there are only two possible envi-

ronments, in which case the agent’s beliefs are modeled by a 445  
mixture of the two. Our exploration policies are independent 446  
of the particular form or implementation of the belief distri- 447  
bution, since they only require the ability to sample from it. 448

We also compare our methods to the model-free minimax 449  
Q-learning algorithm (Littman 1994). The results, in the ap- 450  
pendix, show that our model-based methods are one or two 451  
orders of magnitude more sample efficient! 452

## 453 Conclusions and future research

454 We investigated the increasingly important and common 454  
game-solving setting where we do not have an explicit de- 455  
scription of the game but only oracle access to it through 456  
game play. During a limited-duration learning phase, the al- 457  
gorithm can control the actions of both players in order to 458  
try to learn the game and how to play it well. After that, the 459  
algorithm has to produce a policy that has low exploitabil- 460  
ity. We generalized exploration policies that are used in the 461  
single-agent settings to games. We showed that games raise 462  
additional issues and some generalizations are inherently in- 463  
consistent, and we introduced new ones that fix that. We pro- 464  
posed using the distribution of state-action value functions 465  
induced by a belief distribution over possible environments. 466

467 We conducted experiments on normal-form and stochas- 467  
tic games where exploitability can be evaluated exactly. 468  
The more sophisticated exploration policies—particularly 469  
our Thompson sampling-based policy, which is also com- 470  
putationally efficient—tended to outperform the more naive 471  
exploration policies in solution quality. We also showed one 472  
to two orders of magnitude sample-efficiency improvement 473  
over a well-known model-free exploration approach. 474

475 There are several ways to extend this work. First, one 475  
could replace exact game and game-ensemble solving with 476  
approximate solving, speeding up the computation time of 477  
the exploration algorithms that rely on these solutions. 478

479 Second, one could extend this work to imperfect- 479  
information extensive-form games. For example, one could 480  
use the *sequence form* (von Stengel 1996) representation and 481  
adapt Equation 4 to work over that setting, which is still a 482  
linear program. Alternatively, one could work directly with 483  
the tree form by creating a root chance node whose children 484  
are the games in the ensemble, distinguishing the informa- 485  
tion sets (across the ensemble) that belong to player  $-i$  but 486  
not player  $i$ , and then running a standard game-solving al- 487  
gorithm such as CFR (Zinkevich et al. 2007) or its modern 488  
variants (Brown and Sandholm 2015, 2017, 2019; Brown, 489  
Kroer, and Sandholm 2017; Farina, Kroer, and Sandholm 490  
2021). This will retrieve the argmax-sum-min policy for  $i$ . 491  
Since those algorithms are iterative, one can get fast approxi- 492  
mate solutions (as mentioned earlier) by terminating the pro- 493  
cess early, at the start of each exploration episode. Solutions 494  
from previous episodes could also be used to warm-start the 495  
process (Brown and Sandholm 2016). 496

497 Finally, one could extend this work to environments with 497  
large and/or continuous state spaces where function approx- 498  
imation is required. Our exploration policies suggest 499  
maintaining an ensemble of neural networks modelling the 500  
agent’s belief distribution over environments. Each network 501  
models the transition and reward functions of a hallucinated 502

503 environment, and is repeatedly trained on observation tuples  
504 accumulated over the course of exploration. In parallel, one  
505 trains separate Q-value networks on the inputs and outputs  
506 of each of these environment networks. This yields an en-  
507 semble of Q functions (one for each environment in the be-  
508 lief distribution) that can be used by our exploration algo-  
509 rithms.

## 510 Acknowledgements

511 This material is based on work supported by the Na-  
512 tional Science Foundation under grants IIS-1718457, IIS-  
513 1901403, and CCF-1733556, and the ARO under award  
514 W911NF2010081.

## 515 References

516 Agrawal, S.; and Jia, R. 2017. Optimistic posterior sam-  
517 pling for reinforcement learning: worst-case regret bounds.  
518 In *NIPS*.

519 Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time  
520 analysis of the multiarmed bandit problem. *Machine learn-*  
521 *ing*.

522 Bellemare, M.; Dabney, W.; and Munos, R. 2017. A distri-  
523 butional perspective on reinforcement learning. In *ICML*.

524 Berner, C.; Brockman, G.; Chan, B.; Cheung, V.; Debiak,  
525 P.; Dennison, C.; Farhi, D.; Fischer, Q.; Hashme, S.; Hesse,  
526 C.; Józefowicz, R.; Gray, S.; Olsson, C.; Pachocki, J.; Petrov,  
527 M.; de Oliveira Pinto, H. P.; Raiman, J.; Salimans, T.; Schlatter,  
528 J.; Schneider, J.; Sidor, S.; Sutskever, I.; Tang, J.; Wolski,  
529 F.; and Zhang, S. 2019. Dota 2 with Large Scale Deep Rein-  
530 forcement Learning. *CoRR* URL <http://arxiv.org/abs/1912.06680>.

531 Brown, N.; Kroer, C.; and Sandholm, T. 2017. Dynamic  
532 Thresholding and Pruning for Regret Minimization. In  
533 *AAAI*.

534 Brown, N.; and Sandholm, T. 2015. Regret-Based Pruning  
535 in Extensive-Form Games. In *NIPS*.

536 Brown, N.; and Sandholm, T. 2016. Strategy-Based Warm  
537 Starting for Regret Minimization in Games. In *AAAI*.

538 Brown, N.; and Sandholm, T. 2017. Reduced Space and  
539 Faster Convergence in Imperfect-Information Games via  
540 Pruning. In *ICML*.

541 Brown, N.; and Sandholm, T. 2019. Solving imperfect-  
542 information games via discounted regret minimization. In  
543 *AAAI*.

544 Bubeck, S.; Munos, R.; and Stoltz, G. 2009. Pure Explora-  
545 tion in Multi-armed Bandits Problems. In *ALT*.

546 Casgrain, P.; Ning, B.; and Jaimungal, S. 2019. Deep Q-  
547 Learning for Nash Equilibria: Nash-DQN. *CoRR* URL <http://arxiv.org/abs/1904.10554>.

548 Chen, R.; Sidor, S.; Abbeel, P.; and Schulman, J. 2017. UCB  
549 and InfoGain Exploration via Q-Ensembles. *CoRR* URL  
550 <http://arxiv.org/abs/1706.01502>.

551 Claus, C.; and Boutilier, C. 1998. The dynamics of rein-  
552 forcement learning in cooperative multiagent systems. In  
553 *AAAI*.

554 Conitzer, V.; and Sandholm, T. 2003. BL-WoLF: A Frame-  
555 work For Loss-Bounded Learnability In Zero-Sum Games.  
556 In *ICML*.

557 Dearden, R.; Friedman, N.; and Russell, S. 1998. Bayesian  
558 Q-learning. In *AAAI*.

559 Farina, G.; Kroer, C.; and Sandholm, T. 2021. Faster Game  
560 Solving via Predictive Blackwell Approachability: Connect-  
561 ing Regret Matching and Mirror Descent. In *AAAI*.

562 Ganzfried, S.; and Sandholm, T. 2009. Computing equilibria  
563 in multiplayer stochastic games of imperfect information. In  
564 *IJCAI*.

565 Garivier, A.; Kaufmann, E.; and Koolen, W. M. 2016. Max-  
566 imin Action Identification: A New Bandit Framework for  
567 Games. In *JMLR*.

568 Heinrich, J.; and Silver, D. 2016. Deep Reinforcement  
569 Learning from Self-Play in Imperfect-Information Games.  
570 *CoRR* URL <http://arxiv.org/abs/1603.01121>.

571 Hu, J.; and Wellman, M. 2003. Nash Q-Learning for  
572 general-sum stochastic games. *JMLR*.

573 Kaufmann, E.; Cappe, O.; and Garivier, A. 2012. On  
574 Bayesian upper confidence bounds for bandit problems. In  
575 *AISTATS*.

576 Lanctot, M.; Zambaldi, V.; Gruslys, A.; Lazaridou, A.;  
577 Tuyls, K.; Perolat, J.; Silver, D.; and Graepel, T. 2017. A  
578 Unified Game-Theoretic Approach to Multiagent Reinforce-  
579 ment Learning. In *NIPS*.

580 Littman, M. 1994. Markov games as a framework for multi-  
581 agent reinforcement learning. In *ICML*.

582 Lockhart, E.; Lanctot, M.; Pérolat, J.; Lespiau, J.-B.; Mor-  
583 rill, D.; Timbers, F.; and Tuyls, K. 2019. Computing Ap-  
584 proximate Equilibria in Sequential Adversarial Games by  
585 Exploitability Descent. In *IJCAI*.

586 Marchesi, A.; Trovò, F.; and Gatti, N. 2019. Learn-  
587 ing Probably Approximately Correct Maximin Strategies in  
588 Simulation-Based Games with Infinite Strategy Spaces. In  
589 *AAMAS*.

590 Mavrin, B.; Yao, H.; Kong, L.; Wu, K.; and Yu, Y. 2019. Dis-  
591 tributional reinforcement learning for efficient exploration.  
592 In *ICML*.

593 Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness,  
594 J.; Bellemare, M. G.; Graves, A.; Hiedmiller, M.; Fiedland,  
595 A. K.; Ostrovski, G.; Petersen, S.; Beattie, C.; Sadik, A.;  
596 Antonoglou, I.; King, H.; Kumaran, D.; Wierstra, D.; Legg,  
597 S.; and Hassabis, D. 2015. Human-level control through  
598 deep reinforcement learning. *Nature*.

599 O’Donoghue, B.; Osband, I.; Munos, R.; and Mnih, V.  
600 2018. The uncertainty Bellman equation and exploration.  
601 In *ICML*.

602 Ortega, P. A.; and Braun, D. A. 2010. A minimum relative  
603 entropy principle for learning and acting. *JAIR*.

604 Osband, I.; Blundell, C.; Pritzel, A.; and Roy, B. V. 2016.  
605 Deep exploration via bootstrapped DQN. In *NIPS*.

608 Osband, I.; and Roy, B. V. 2017. Why is Posterior Sampling Better than Optimism for Reinforcement Learning? In *PMLR*.

609

610

611 Osband, I.; Roy, B. V.; Russo, D. J.; and Wen, Z. 2019. Deep exploration via randomized value functions. *JMLR*.

612

613 Russo, D.; van Roy, B.; Kazerouni, A.; Osband, I.; and Wen, Z. 2018. A Tutorial on Thompson Sampling. *Foundations and Trends in Machine Learning*.

614

615

616 Sandholm, T.; and Crites, R. 1996. Multiagent Reinforcement Learning in the Iterated Prisoner’s Dilemma. *Biosystems*.

617

618

619 Shapley, L. 1953. Stochastic Games. *PNAS*.

620

621 Sokota, S.; Ho, C.; and Wiedenbeck, B. 2019. Learning deviation payoffs in simulation-based games. In *AAAI*.

622

623 Srinivasan, S.; Lanctot, M.; Zambaldi, V.; Perolat, J.; Tuyls, K.; Munos, R.; and Bowling, M. 2018. Actor-Critic Policy Optimization in Partially Observable Multiagent Environments. In *NeurIPS*.

624

625

626 Thompson, W. 1933. On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika*.

627

628

629 Vinyals, O.; Babuschkin, I.; Czarnecki, W. M.; Mathieu, M.; Dudzik, A.; Chung, J.; Choi, D. H.; Powell, R.; Ewalds, T.; Georgiev, P.; Oh, J.; Horgan, D.; Kroiss, M.; Danihelka, I.; Huang, A.; Sifre, L.; Cai, T.; Agapiou, J. P.; Jaderberg, M.; Vezhnevets, A. S.; Leblond, R.; Pohlen, T.; Dalibard, V.; Budden, D.; Sulsky, Y.; Molloy, J.; Paine, T. L.; Gulcehre, C.; Wang, Z.; Pfaff, T.; Wu, Y.; Ring, R.; Yogatama, D.; Wunsch, D.; McKinney, K.; Smith, O.; Schaul, T.; Lillicrap, T.; Kavukcuoglu, K.; Hassabis, D.; Apps, C.; and Silver, D. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*.

636

637

638

639

640 von Stengel, B. 1996. Efficient Computation of Behavior Strategies. *Games and Economic Behavior*.

641

642

643 Vorobeychik, Y.; and Wellman, M. 2009. Strategic analysis with simulation-based games. In *WSC*.

644

645

646

647 Wang, X.; and Sandholm, T. 2002. Reinforcement Learning to Play An Optimal Nash Equilibrium in Team Markov Games. In *NIPS*.

648

649

650

651

652

653

654

655

656

657

## Appendix

In this appendix we present additional technical material that did not fit in the body of the paper.

## Upper bound on the suboptimality of the argmax-sum-min estimator

658  
659

**Theorem 1.** Let  $f : X \times Y \rightarrow \mathbb{R}$ . Let  $D$  be a probability distribution on  $Y$ . Let

$$\hat{x} \in \operatorname{argmax}_x \frac{1}{n} \sum_i f(x, y_i) \quad (27)$$

where  $i \in [n], y_i \sim D$ . Let  $y \sim D$ . The probability

$$\max_x \mathbb{E}_y f(x, y) - \mathbb{E}_y f(\hat{x}, y) \quad (28)$$

exceeds  $\varepsilon$  is at most  $e^{-2n\varepsilon^2/\Delta^2}$ , where  $\Delta$  is the range (supremum minus infimum) of  $f$ .

660  
661

*Proof.* By Hoeffding’s inequality,

$$p = \Pr_{y_i} \left( \mathbb{E}_y f(x, y) - \mathbb{E}_y f(\hat{x}, y) \geq \varepsilon \right) \quad (29)$$

$$= \Pr_{y_i} \left( \mathbb{E}_y f(x, y) - \mathbb{E}_y f(\hat{x}, y) + \frac{1}{n} \sum_i f(x, y_i) \right) \quad (30)$$

$$- \frac{1}{n} \sum_i f(x, y_i) \geq \varepsilon \right) \quad (31)$$

$$\leq \Pr_{y_i} \left( \mathbb{E}_y f(x, y) - \mathbb{E}_y f(\hat{x}, y) + \frac{1}{n} \sum_i f(\hat{x}, y_i) \right) \quad (32)$$

$$- \frac{1}{n} \sum_i f(x, y_i) \geq \varepsilon \right) \quad (33)$$

$$= \Pr_{y_i} \left( \mathbb{E}_y (f(x, y) - f(\hat{x}, y)) \right) \quad (34)$$

$$- \frac{1}{n} \sum_i (f(x, y_i) - f(\hat{x}, y_i)) \geq \varepsilon \right) \quad (35)$$

$$\leq \exp \left( - \frac{2n\varepsilon^2}{\Delta^2} \right) \quad (36)$$

□ 662

## Additional figures

663

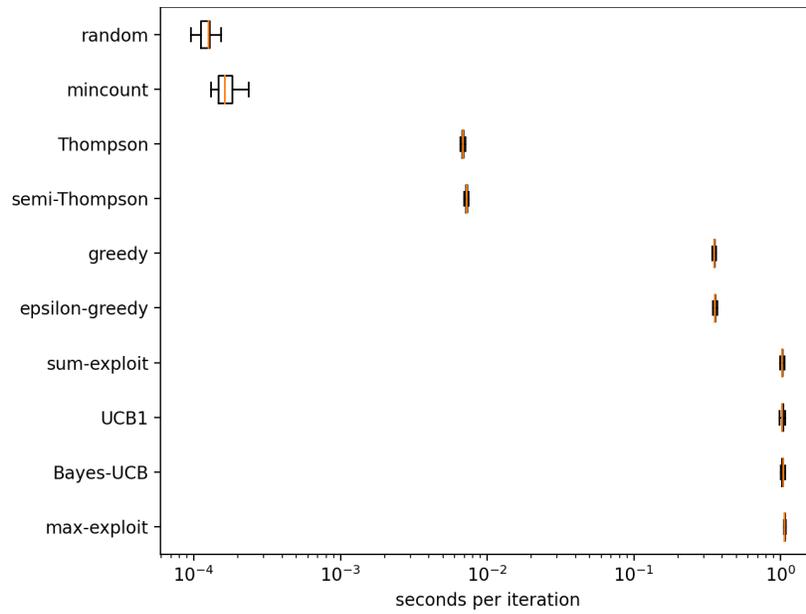


Figure 2: Typical time to compute each exploration policy at the beginning of each exploration episode on a randomly-generated  $10 \times 10$  normal-form game.

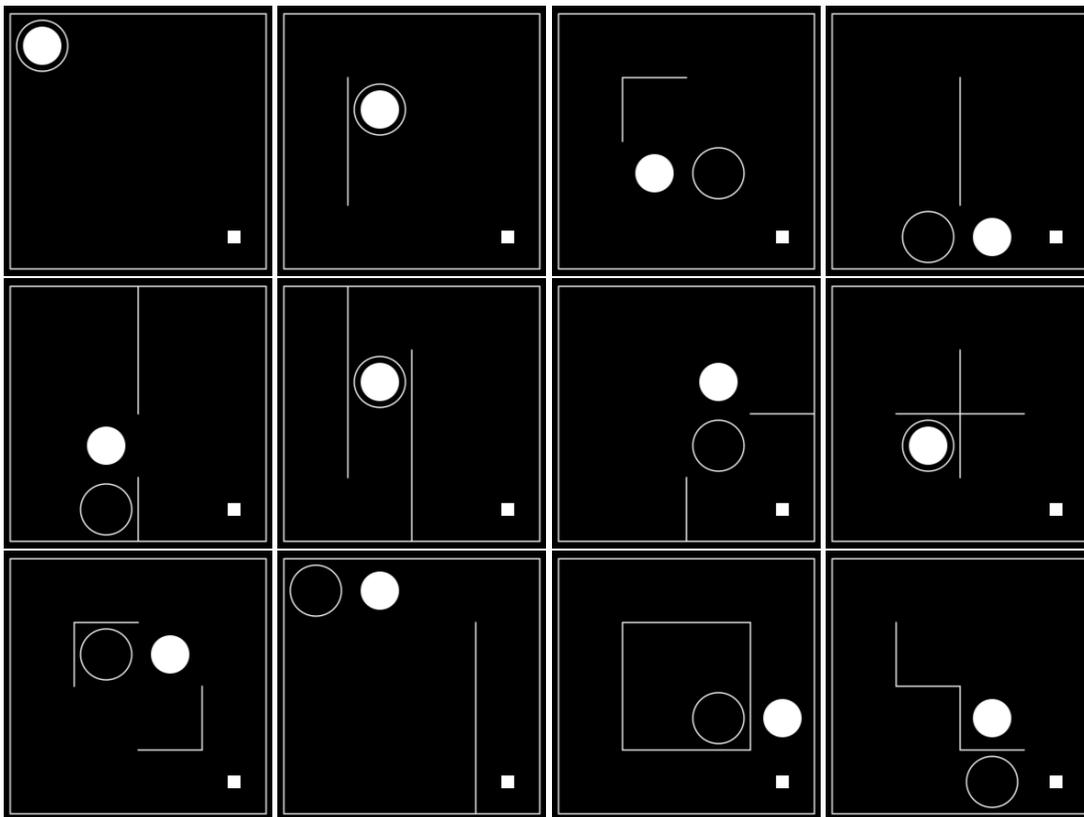


Figure 3: Layouts used for the bombardment environment. Left-to-right, top-to-bottom: empty, corridor, fork, wall, gap, snake, nest, cross, double fork, strip, box, stair.

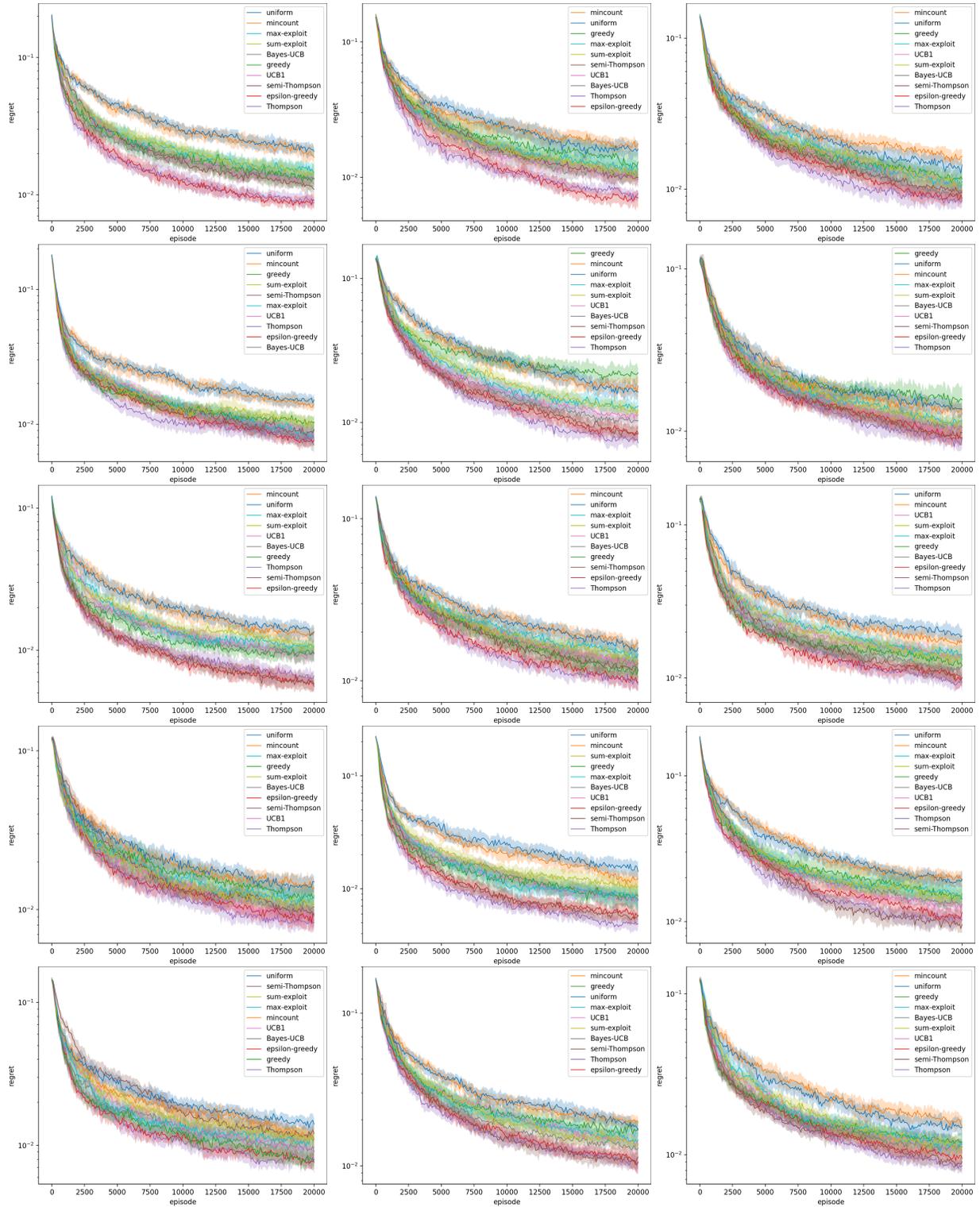


Figure 4: Performance of each policy on more randomly-generated  $10 \times 10$  normal-form games. Rewards are Bernoulli random variables whose bias is sampled from the uniform distribution, which is also the prior used by the exploratory agent. 100 trials were used. For the policies requiring it, 100 game samples and 100 strategy or strategy profile samples were used on each iteration. For the  $\epsilon$ -greedy we used  $\epsilon = 0.1$ . The solid curves show the medians and the shaded regions show the 40-60 percentile bands.

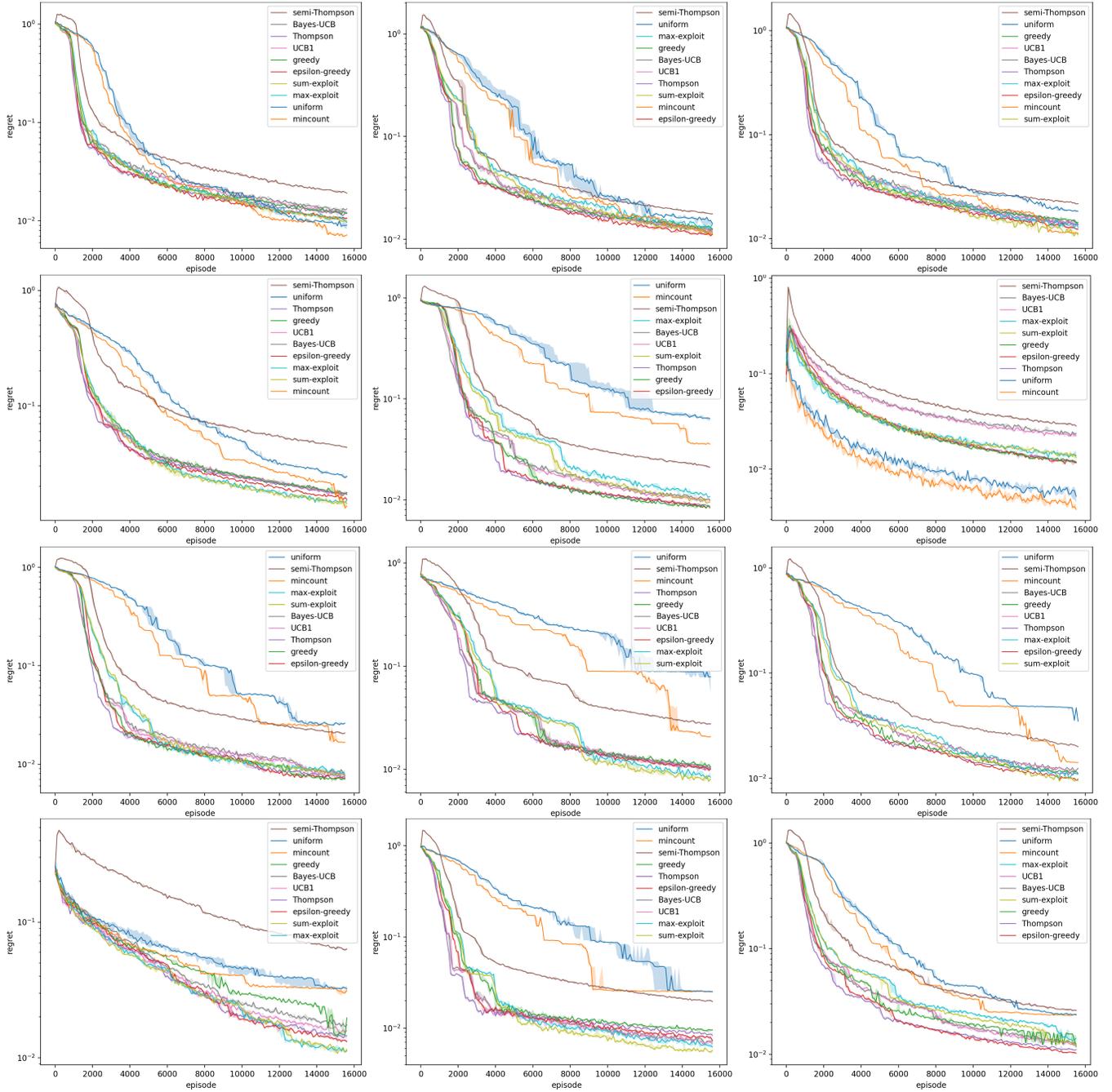


Figure 5: Performance of each policy on each bombardment game layout. Left-to-right, top-to-bottom: empty, corridor, fork, wall, gap, snake, nest, cross, double fork, strip, box, and stair. Each episode had a 40-step horizon. The solid curves show the medians and the shaded regions show the 25-75 percentile bands. Other parameters are the same as for the normal-form game setting.

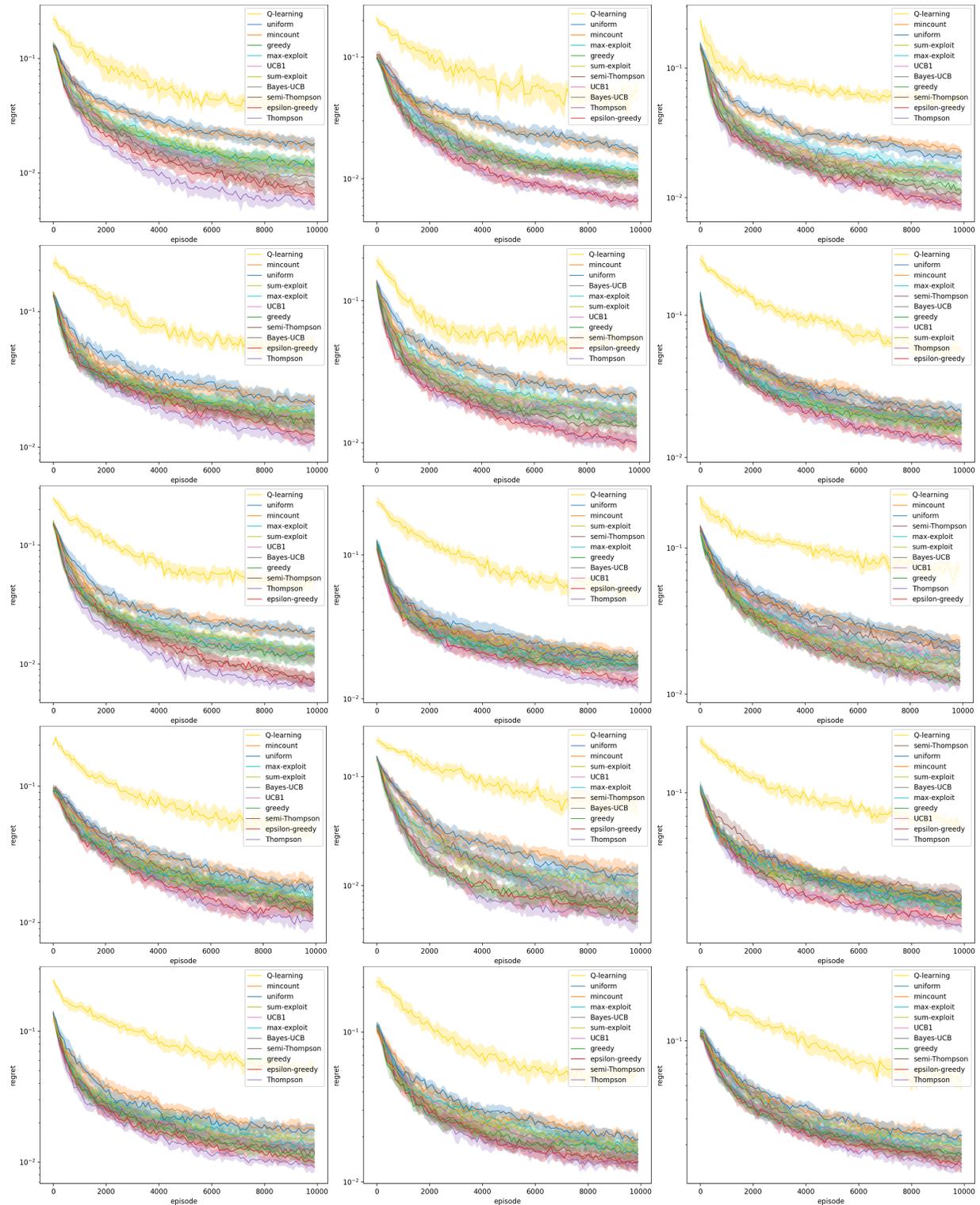


Figure 6: Performance of minimax Q-learning with learning rate 0.1 and  $\epsilon$ -greedy exploration ( $\epsilon = 0.1$ ) on randomly-generated normal form games. Other parameters are the same as for Figure 4.

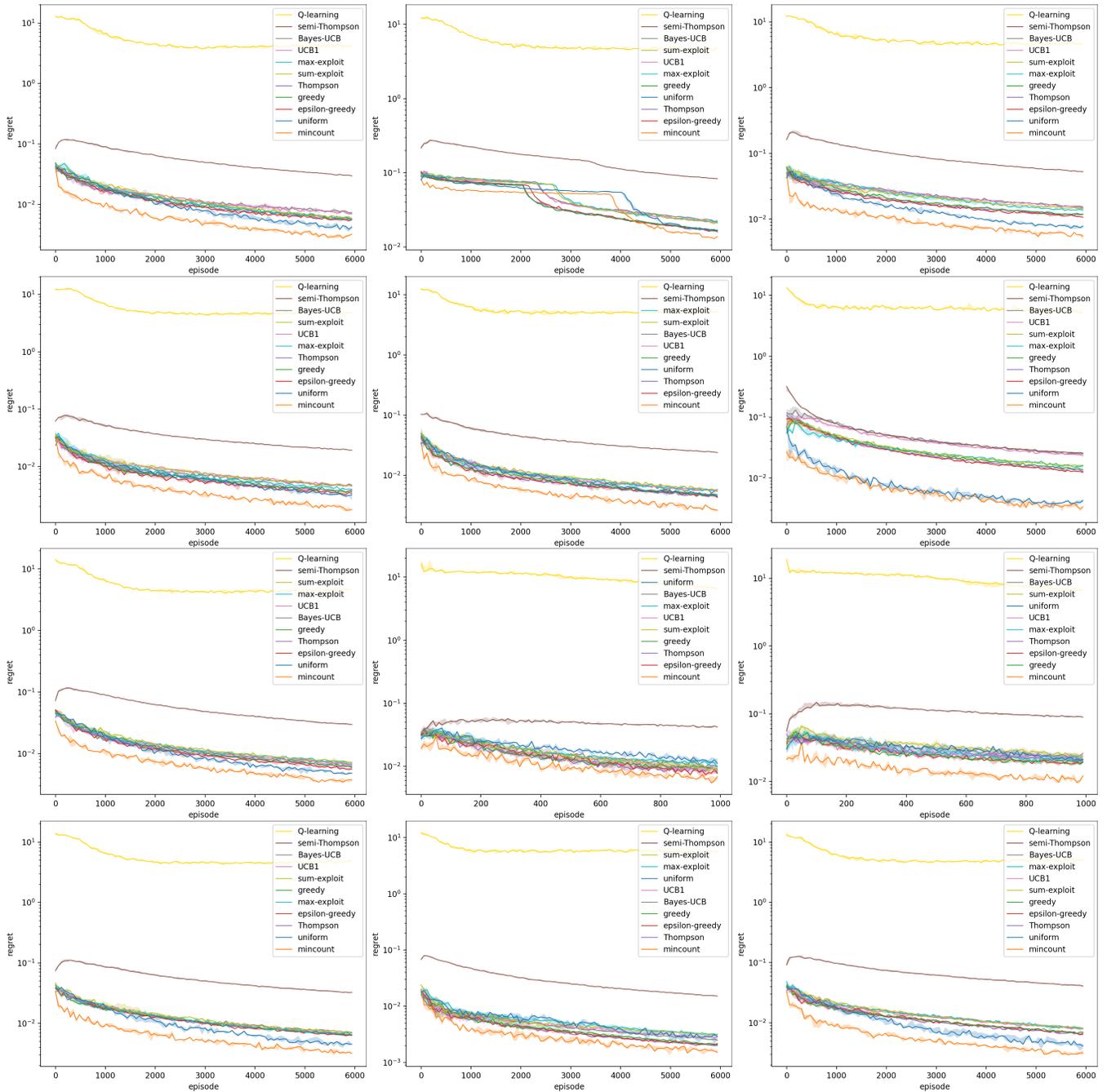


Figure 7: Performance of minimax Q-learning with learning rate 0.1 and  $\epsilon$ -greedy exploration ( $\epsilon = 0.1$ ) on each bombardment game. Other parameters are the same as for Figure 5.