# **Composing Efficient, Robust Tests for Policy Selection**

Dustin Morrill<sup>1</sup> Thomas J. Walsh<sup>1</sup> Daniel Hernandez<sup>1</sup> Peter R. Wurman<sup>1</sup> Peter Stone<sup>1,2</sup>

<sup>1</sup>Sony AI, New York, NY, USA <sup>2</sup>Department of Computer Science, The University of Texas at Austin, Austin, TX USA

### **Abstract**

Modern reinforcement learning systems produce many high-quality policies throughout the learning process. However, to choose which policy to actually deploy in the real world, they must be tested under an intractable number of environmental conditions. We introduce RPOSST, an algorithm to select a small set of test cases from a larger pool based on a relatively small number of sample evaluations. RPOSST treats the test case selection problem as a two-player game and optimizes a solution with provable k-of-N robustness, bounding the error relative to a test that used all the test cases in the pool. Empirical results demonstrate that RPOSST finds a small set of test cases that identify high quality policies in a toy one-shot game, poker datasets, and a high-fidelity racing simulator.

## 1 INTRODUCTION

Reinforcement learning (RL) [Sutton and Barto, 2018] policies have made a number of stunning breakthroughs in multiplayer games [Silver et al., 2016, Moravčík et al., 2017, Brown and Sandholm, 2018, Vinyals et al., 2019, Brown and Sandholm, 2019, Wurman et al., 2022, FAIR et al., 2022, Perolat et al., 2022]. However, the process of choosing an RL policy for production usage, either in an exhibition or deployment for end users, is challenging. Practitioners often generate many policies that perform well during training but which require thorough vetting on alternative conditions or opponents. Ideally, we would construct a test case for every conceivable deployment scenario, evaluate each policy on each test case, and rank each policy according to a weighted average of test case results. However, such a procedure is typically infeasible because of the sheer numbers of policies and deployment scenarios, especially if test cases are lengthy or involve people. In this work, we present a method

for selecting a small number of test cases from a larger pool that minimizes the reduction in test quality.

Practitioners from other fields, e.g., educational testing [van der Linden, 2005], will recognize this problem as test construction—selecting a small yet robust set of test cases, based on limited data, to evaluate many candidates. This set of test cases should contain enough information to indicate performance over the whole test case pool. For instance, if a policy can defeat a skilled opponent, we can infer that it can defeat an unskilled opponent. However, complicated domains contain complex intransitive relationships between policies, necessitating test case diversity. In addition, there is considerable uncertainty over what policies may be produced in the future and what test cases are the most important to game designers. This uncertainty needs to be considered because once test cases are chosen, the future policies to assess may be the most difficult ones for the test to evaluate accurately. Therefore, a robust solution is required.

We introduce a framework, *robust population optimization for a small set of test cases* (*RPOSST*), to compose an efficient robust test of a fixed size. RPOSST tunes its test to approximate the test scores of adversarially selected policies and test case averaging weights, given test case results on a small set of policies. We present two RPOSST algorithms representing different use cases, focusing on RPOSST<sub>SEQ</sub>, which is better suited to current RL deployment pipelines. We provide robustness guarantees for RPOSST<sub>SEQ</sub> and CVaR RPOSST<sub>SEQ</sub> (a convenient special case) for *k*-of-*N* robustness measures [Chen and Bowling, 2012]. These guarantees provide confidence that RPOSST test scores for future deployment candidates are reliable.

Our contributions include the RPOSST framework, including two algorithm versions, robustness guarantees, and empirical validation in domains widely ranging in complexity. Empirical results are presented for a toy one-shot game simulating race car passing, computer poker competition datasets, and the high fidelity racing simulator, Gran Turismo<sup>TM</sup> 7.

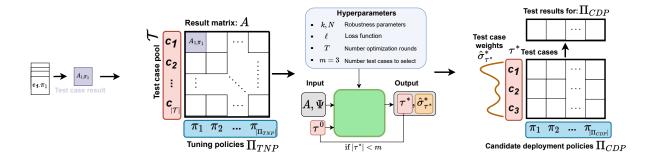


Figure 1: Policy testing with RPOSST. From left to right, the result matrix A is constructed from rollouts, i.e.,  $A_{i,j}$  is the average rollout outcome for tuning policy j ( $\Pi_{\text{TNP}}$  is the set of tuning policies) on test case i. RPOSST analyzes A, taking into account uncertainty distribution  $\Psi$  and a (possibly empty) initial set of test cases that must be used,  $\tau^0$ . RPOSST outputs an efficient robust test  $\langle \tau^*, \hat{\sigma}_{\tau^*}^* \rangle$ , here using only m=3 test cases (if  $\mathcal{T}$  is too large to select all 3 test cases at once,  $\tau^*$  can be fed back into RPOSST as  $\tau^0$ ). New candidate deployment policies,  $\Pi_{\text{CDP}}$ , are tested against each test case in  $\tau^*$  and each result is weighted according to  $\hat{\sigma}_{\tau^*}^*$ , producing a test score for each candidate deployment policy.

They show that RPOSST can dramatically reduce (compared to the full set) the number of test cases needed to identify good deployment policies.

## 2 PROBLEM DEFINITION

The goal of *policy testing* is to evaluate the strengths and weaknesses of a large set of *candidate deployment policies*,  $\Pi_{\text{CDP}} \subseteq \Pi$ , in order to choose one for deployment. A *policy*  $\pi \in \Pi$  in this setting can be any mapping from environment observations to a distribution over actions (*e.g.*, Markov policies; Sutton and Barto [2018]). A policy is evaluated on a *test* consisting of various test cases chosen from a pool,  $\mathcal{T}$ . Each *test case* simulates an important aspect of the deployment environment, for example, different parameter settings like weather conditions or different opponent policies in a competative game. For straightforward comparisons between policies, we summarize a policy  $\pi$ 's test results with a scalar *test score*, computed as the weighted average of  $\pi$ 's test case results according to *test case weights*,  $\sigma \in \triangle^{|\mathcal{T}|}$ .

If  $\mathcal{T}$  is small, then right before deployment we could simply test each policy, rank the policies in  $\Pi_{\text{CDP}}$  according to the test scores, and deploy the best one. However, if policies will encounter a wide range of conditions during deployment, e.g., hundreds or thousands of different players for a policy deployed to a popular video game, then  $\mathcal{T}$  ostensibly needs to be large in order to adequately reflect such diversity. The linear scaling in  $|\mathcal{T}|$  presents not just a computational burden, but also costs in sample complexity (if the test cases are lengthy) or even in person-time if human quality assurance testers might be needed for test cases.

This work addresses the problem of composing an efficient test,  $\langle \tau, \hat{\sigma}_{\tau} \rangle$ , by selecting a small number of test cases  $\tau \subset \mathcal{T}$  and test case weights  $\hat{\sigma} \in \triangle^{|\tau|}$  to approximate a full test,  $\langle \mathcal{T}, \sigma \in \triangle^{|\mathcal{T}|} \rangle$ . Complicating this task are two sources of uncertainty to which the efficient test must be robust. First,  $\langle \tau, \hat{\sigma}_{\tau} \rangle$  ought to be used on new candidate deployment policies, so  $\Pi_{\text{CDP}}$  is unknown before  $\langle \tau, \hat{\sigma}_{\tau} \rangle$  is chosen. Second, the desired target ta

We assume access to a small set of representative *tuning* policies  $\Pi_{\text{TNP}} \subset \Pi$  for immediate testing (Section 4 discusses practical considerations in the composition of  $\Pi_{\text{TNP}}$ ). Additionally, our algorithm takes as input a joint distribution  $\Psi$  over  $\Pi_{\text{TNP}}$  and  $\triangle^{|\mathcal{T}|}$  to represent the combined uncertainty about which policies the output test will be applied to and which target distribution to approximate. See Figure 1 for an illustration of the test composition pipeline.

As a concrete example of the terms above and the need for robustness in the face of uncertainty, consider a car-racing agent developed for a one-on-one racing game. The first source of uncertainty is over the future policies we may want to test. Consider the case where, at test construction time, we have policies from two training runs—one that produces aggressive (collision-prone) policies, and another that produces more polite policies, but we are uncertain about which type will be best suited for the game. In this case, we want the selected test cases to provide good evaluations on policies from either set, and thus require  $\Psi$  to reflect this uncertainty. Policies from both sets should be included in  $\Pi_{\text{TNP}}$  and our algorithm needs to be robust to policies within  $\Pi_{\text{TNP}}$ .

The second source of uncertainty is over which test cases are most important. Imagine that we have some test cases that specifically target and penalize off-track infractions. In the future, game designers could request fewer infractions or allow for more risky racing lines. To hedge against both of these possibilities we can add two target distributions to  $\Psi$ , one where off-track tests cases have higher weights than the other test cases and another where they have lower weights. The job of an algorithm (such as RPOSST) is then to ensure its tests are accurate according to both target distributions.

## 3 BACKGROUND

In order to compose an efficient and robust test, we utilize established game-theoretic frameworks for modeling robustness and learning optimal decisions (specifically, regret minimization). The following subsections present background material on these two topics.

#### 3.1 ROBUSTNESS

The idea of *robustness* is to prepare for an unfavorable portion of possible outcomes sampled from an uncertainty distribution. In our formulation of policy testing the uncertainty distribution covers the future policies in  $\Pi_{\text{CDP}}$  and the target distribution. A *percentile robustness measure* [Charnes and Cooper, 1959],  $\mu$ , is a formal representation of a robustness criterion as a probability distribution over percentiles. For example, if  $\mu$  has all of its weight on 0.01, then an m-size test with weights  $\hat{\sigma}_{\tau}$  that is robust according to  $\mu$ , then the test minimizes test score error on  $\hat{\sigma}_{\tau}$ 's worst 1% of policy–target-distribution pairs sampled from  $\Psi$ .

The k-of-N robustness measures [Chen and Bowling, 2012] are percentile robustness measures defined by parameters  $k,N\in\mathbb{N},\,1\leq k\leq N$ , that permit tractable optimization procedures. This parameterization reflects the mechanics of how an efficient test  $\langle \tau,\hat{\sigma}_{\tau}\rangle$  is evaluated on such a measure: N policy-target-distribution pairs are sampled from  $\Psi$  and  $\hat{\sigma}_{\tau}$ 's performance is averaged over the k worst pairs for  $\hat{\sigma}_{\tau}$ . Every k-of-N robustness measure is a non-increasing function, i.e., more weight is placed on smaller percentiles, and the fraction k/N represents the percentile (technically the fractile) around which the measure decreases.

In our test construction setting, the choice of k and N reflects the designer's tolerance for test scores that are bad because of "unlucky" outcomes from  $\Psi$  (that is, test scores with large error on policy–target-distribution pairs sampled from  $\Psi$ , even if they are sampled infrequently). Optimizing for performance under small percentiles (e.g., setting  $k=1,\ N=100$ ) yields tests with a small maximum test score error across  $\Pi_{\rm TNP}$ . Then, even if each candidate deployment policy resembles the tuning policy that has the largest test score error, the optimized test will yield small

test score errors. In contrast, optimizing for the uniform measure (k=N) optimizes for mean performance across  $\Pi_{\text{TNP}}$ , essentially assuming  $\Pi_{\text{CDP}}=\Pi_{\text{TNP}}$ , which can lead to large test score error on the actual candidate deployment policies.

As  $N \to \infty$ , the k-of-N robustness measure approaches the *conditional value at risk* (CVaR) robustness measure at the k/N fractile [Chen and Bowling, 2012], which evenly weights all of the fractiles  $\leq k/N$  and puts a weight of zero on all larger fractiles. Formally, the robustness optimization objective is to minimize the *percentile performance loss*:

$$L_{\mu,\Psi}(\hat{\sigma}_{\tau}) = \inf_{y \in \mathcal{Y}} \int_{\eta \in [0,1], \mathbb{P}[\ell(\hat{\sigma}_{\tau};\pi,\sigma) \le y(\eta)] \ge \eta} y(\eta)\mu(d\eta), \tag{1}$$

under a loss function  $\ell: \triangle^{|\mathcal{T}|} \times \Pi_{\mathsf{TNP}} \times \triangle^{|\mathcal{T}|} \to \mathbb{R}$  where we overload  $\ell$  for incomplete test case weight vectors by filling in zeros for missing elements,  $\langle \pi, \sigma \rangle \sim \Psi$ , and  $\mathcal{Y}$  is the class of real-valued, bounded,  $\mu$ -integrable functions on [0,1]. An *efficient* (m-size)  $\mu$ -robust test is a minimizer of  $L_{\mu,\Psi}$  across all  $\hat{\sigma}_{\tau}$  where  $\tau=m$ .

The optimization of the percentile performance loss under k-of-N robustness measure,  $\mu_{k$ -of-N, can be modeled as a zero-sum imperfect information game [Chen and Bowling, 2012]. Here, a protagonist player constructs efficient tests and an antagonist chooses a tuning policy to test and a target distribution. For their payoffs, the antagonist receives the test score error of the protagonist's test given the antagonist's tuning policy and target distribution while the protagonist receives the negation. The k and N parameters determine which target distributions and tuning policies that the antagonist can choose from and how many pairs must be averaged across. At the start of the game, N targetdistribution—tuning-policy pairs are sampled. From these Npairs, the antagonist must select k of them. Finally, one of these k pairs is sampled, both players receive their payoffs, and the game ends. A *minimax* test for the protagonist, *i.e.*, one that minimizes the protagonist's maximum loss in this game is a  $\mu_{k\text{-of-}N}$ -robust test.

#### 3.2 REGRET

While the game above models the optimization process, it does not instruct the protagonist on *how* to choose test cases to win. A no-regret *online decision process* (*ODP*) algorithm can find approximate minimax decisions by repeatedly playing out the game and improving over time from payoff feedback. Formally, on each round t of the game, an ODP algorithm chooses an efficient test  $\langle \tau^t, \hat{\sigma}_{\tau^t}^t \rangle$  and receives the *payoff function*  $v^t = -\nabla_{\hat{\sigma}_{\tau^t}^t} \ell(\hat{\sigma}_{\tau^t}^t; \pi^t, \sigma^t)$  as feedback given  $\langle \pi^t, \sigma^t \rangle$  chosen by the antagonist. If the antagonist always plays a best response to the ODP algorithm, that is, the tuning-policy-target-distribution pair that maximizes the loss of  $\hat{\sigma}_{\tau^t}^t$  on each round  $t \in \{1, \dots, T\}, T \geq 1$ , then

the *no-regret* property ensures that at least one of the tests in the sequence  $\langle\langle \tau^t, \hat{\sigma}^t_{\tau^t} \rangle\rangle_{t=1}^T$  is at most  $\mathcal{O}(G/\sqrt{T})$  away from the minimax value, where G>0 is the maximum magnitude of the loss gradient (see Lockhart et al. [2019a,b] and Appendix Proposition C.4 for more details).

Regret matching<sup>+</sup> [Tammelin, 2014, Tammelin et al., 2015] is a no-regret algorithm for simplex decision sets, e.g., the m dimensional test case weight space  $\triangle^m$ , that selects  $\hat{\sigma}^t_{\tau^t} = \frac{q^{1:t-1}}{\mathbf{1}^+ q^{1:t-1}}$  using pseudoregrets  $q^{1:t} = [q^{1:t-1} + \rho^t]_+, q^{1:0} = \mathbf{0}$ , where  $\rho^t = v^t - (v^t)^\top \hat{\sigma}^t_{\tau^t}$  is the instantaneous regret vector ( $\hat{\sigma}^t_{\tau^t} = \frac{1}{d}\mathbf{1}$  if none of the pseudoregrets are positive).

## 4 RPOSST

Our approach, robust population optimization for a small set of test cases (RPOSST) begins by evaluating each tuning policy  $\pi \in \Pi_{\mathsf{TNP}}$  on each test case  $c \in \mathcal{T}$ , yielding a  $|\mathcal{T}| \times |\Pi_{\mathsf{TNP}}|$  result matrix A of test case results. As an optimization approach, RPOSST aims to minimize prediction errors, as measured by a convex function  $\Delta : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ , e.g., the absolute difference  $\Delta(\hat{x},x) = |\hat{x}-x|$ . RPOSST robustly optimizes for a small set of test cases and a weighting over them according to how well it reproduces test scores admitted by A as measured by a loss function

$$\ell: \hat{\sigma}; \pi_j, \sigma \mapsto \Delta(\underbrace{E_{i \sim \hat{\sigma}}[A_{i,j}]}_{\hat{\sigma}\text{'s test score for } \pi_j.}, \underbrace{E_{i \sim \sigma}[A_{i,j}]}_{\text{o's test score for } \pi_j.}),$$

on test case distribution  $\hat{\sigma} \in \triangle^{|\mathcal{T}|}$  compared to  $\sigma \in \triangle^{|\mathcal{T}|}$  with respect to test results from the  $j^{\text{th}}$  tuning policy  $\pi_j$ . Since  $\hat{\sigma}$  is being used to produce test scores that approximate those under  $\sigma$ , we call  $\sigma$  a target distribution in this context. Our goal is to select a small number of test cases, so we constrain RPOSST to output weights  $\hat{\sigma}_{\tau} \in \triangle^m$  for groups of test cases  $\tau \subset \mathcal{T}$  of size m.

Though  $\mathcal{T}$  is large, the cost of computing A is balanced by the savings of using fewer test cases for future policies. RPOSST is robust to any distribution over  $\Pi_{TNP}$ , so as long as this set covers the space of  $\Pi_{CDP}$  (i.e., all  $\pi \in \Pi_{CDP}$ are convex mixtures of  $\Pi_{TNP}$ ), this robustness imparts a minimum test accuracy guarantee even on deployment candidates. Intuitively, this means the quality of RPOSST's tests will tend to improve with more diverse tuning policies. Accordingly, it should be beneficial for a tuning policy to represent an extreme point in a reasonable region of policy space, or at least for it to be generated with a method similar to that which will generate deployment candidates (e.g., sampled from checkpoints of RL training runs). That way, the tuning policies include a diverse collection of skilled and unskilled policies with random variations, while retaining architectural and algorithmic similarities to future deployment candidates.

Following the earlier discussion of k-of-N robustness, we frame the optimization in RPOSST as a zero-sum game.

By adversarially choosing policies to test, the antagonist forces RPOSST to compose tests that are better at accurately testing the more difficult-to-assess policies in the tuning set, providing a degree of robustness to the distribution of future deployment candidates. Similarly, by adversarially choosing the target distribution, the antagonist also forces RPOSST to be robust along this dimension. The steps of each round  $t=1,\ldots,T$  of our optimization game follows.

- 1. The protagonist must choose an m-tuple of test cases  $\tau^t \subset \mathcal{T}$  and weights  $\hat{\sigma}_{\tau^t}^t \in \triangle^m$ .
- 2. N policies to test and target distributions,  $\langle\langle \pi_{j_i}, \sigma_i \rangle\rangle_{i=1}^N$ , are sampled from uncertainty distribution  $\Psi$ .
- 3. The antagonist chooses the k worst policies and target distributions, *i.e.*, those that maximize  $\ell(\hat{\sigma}_{\tau t}^t; \pi_{j_i}, \sigma_i)$ .
- 4. One of the k worst configurations is sampled uniformly, leading to the end of the round, at which point the protagonist receives the payoff  $v_{\tau^t,(i)}^t = -\ell(\hat{\sigma}_{\tau^t}^t; \pi_{j_{(i)}}, \sigma_{(i)})$ , where the subscript (i) denotes the  $i^{\text{th}}$  element of a sorted list in descending order (the  $i^{\text{th}}$  worst for the protagonist).

The protagonist is allowed to update their strategy at the end of each round based on the expected payoff,  $\mathbb{E}_{i \sim \text{Unif}(\{1,\dots,k\})} \left[ v_{\tau,(i)}^t \right]$ , for each  $\tau \in \mathcal{T}^m$  they could have chosen. The more rounds of the game that are run (the larger T is), the closer RPOSST gets to returning a minimax strategy, and consequently, a robust optimal selection of test cases and weights. Thus, in application, T can be set as large as is convenient under computational and time constraints. Theorem D.2 gives a precise rate for RPOSST's improvement, with high probability, as a function of T. Although the protagonist must consider an exponential (in m) number of test case combinations, the premise of RPOSST is that we want a small set of test cases, so m will be small. To decrease computational requirements, RPOSST can be run in a loop to select test cases iteratively until m have been selected, at a potential cost to test accuracy compared to optimizing for the entire m-tuple at once.

## 4.1 ANTAGONIST INFORMATION MODELS

We consider two RPOSST algorithm variants that utilize different models of the information that the antagonist in our optimization game has before they make their choice. These models correspond to two policy testing use cases. The first, "simultaneous move" model is less pessimistic, but has impractical aspects, which are addressed by the subsequent "sequential move" model.

**Simultaneous move.** The simultaneous move model is a naïve application of the original k-of-N game by [Chen and Bowling, 2012]. In this model, the antagonist does not

```
1 Inputs: \langle k, N, T_1, m, \Psi, \tau^0, \ell, T_2 \rangle
 q_{\tau}^{1:0} \leftarrow \mathbf{0} \in \mathbb{R}^{m+|\tau^0|} \text{ for } \tau \in \mathcal{T}^m
 T' \sim \text{Unif}(\{1, ..., T_1\})
 4 for t \leftarrow 1, \dots, T' do
                  for 	au \in \mathcal{T}^m do
                            \begin{aligned} z^t &\leftarrow \mathbf{1}^\top q_\tau^{1:t-1} \\ \hat{\sigma}_\tau^t &\leftarrow q_\tau^{1:t-1}/z^t \text{ if } z^t > 0 \text{ else } \mathbf{1}/m \end{aligned}
  6
  7
                             // Add zeros to ensure \hat{\sigma}_{\tau}^{'} \in \triangle^{|\mathcal{T}|} .
  8
                             \hat{\sigma}_{\tau}^t(x) \leftarrow 0 \text{ for } x \in \mathcal{T} \setminus (\tau \cup \tau^0)
                          [\ell_{\tau,(i)}]_{i=1}^{k} \leftarrow \mathcal{L}_{\text{k-of-N}}(\hat{\sigma}_{\tau}^{t}, \langle k, N \rangle, \Psi, \ell)
v_{\tau}^{t} \leftarrow \frac{-1}{k} \sum_{i=1}^{k} \frac{\partial \ell_{\tau,(i)}, \pi_{j_{(i)}}}{\partial \hat{\sigma}_{\tau}^{t}}
10
11
                     // Update regret matching<sup>+</sup>.  \begin{aligned} \rho_{\tau}^t &\leftarrow v_{\tau}^t - (\hat{\sigma}_{\tau}^t)^\top v_{\tau}^t \\ q_{\tau}^{1:t} &\leftarrow [q_{\tau}^{1:t-1} + \rho_{\tau}^t]_+ \end{aligned} 
12
          \mathrm{SR}\bigg(\tau \mapsto \frac{1}{2kL}\mathbf{1}^{\top}\,\mathcal{L}_{\text{k-of-N}}\Big(\hat{\sigma}_{\tau}^{T'},\langle k,N\rangle,\Psi,\ell\Big),T_2\bigg)
16 return 	au^*, \hat{\sigma}_{	au^*}^{T'}
 1 Procedure \mathcal{L}_{\text{k-}\textit{of-N}} Inputs: \langle \hat{\sigma}, k, N, \Psi, \ell \rangle
                  for i = 1 \dots N do
  3
                             // Sample antagonist actions.
  4
                            // Evaluate \hat{\sigma}.
  5
                      \ell_i \leftarrow \ell(\hat{\sigma}; \pi_{j_i}, \sigma_i)
  6
                  \ell_i \leftarrow \ell(\sigma; \pi_{j_i}, \sigma_i)
// Sort to identify the worst k.
                 \operatorname{Sort}\left(\left[\ell_{i}\right]_{i=1}^{N}\right)
return \left[\ell_{(i)}\right]_{i=1}^{k}
```

**Algorithm 1:** RPOSST<sub>SEQ</sub> with regret matching<sup>+</sup> and Successive Rejects

observe which m-tuple of test cases,  $\tau^t$ , is selected by the protagonist on each round t. Instead, it is randomized with a distribution  $\hat{\sigma}_{\mathcal{T}}^t \in \triangle^{|\mathcal{T}|^m}$ . This model corresponds to the policy testing use case where a new m-tuple of test cases is sampled independently for each test that is performed. Every test only evaluates m cases, as desired from a computational efficiency perspective, however, the particular test cases used in each test could be different, making results incomparable across tests. See Appendix Appendix F for additional details.

Sequential move. In the sequential move model, the antagonist observes  $\tau^t$  before acting. The antagonist is thus able to tailor their choice of  $\left\langle \left\langle \pi_{j_{(i)}}, \sigma_{(i)} \right\rangle \right\rangle_{i=1}^k$  to whichever  $\tau^t$  is selected, and randomizing over the m-tuple of test cases has no benefit to the protagonist. Since the antagonist observes  $\tau^t$ , the protagonist must update all the weights that they would apply to each test case tuple  $\tau$  as if  $\tau^t = \tau$ . Thus, the selection of  $\tau^t$  does not impact the protagonist's updates and we need not explicitly select an m-tuple until the very end

of the algorithm, after  $T' \sim \text{Unif}(\{1, \dots, T_1\})$  rounds.<sup>1</sup>

Since the set of N losses observed on each round are generally random, we cannot reuse them to identify which mtuple leads to the lowest loss using the the test case weights computed after running for T' rounds,  $\langle \hat{\sigma}_{\tau}^{T'} \rangle_{\tau \in T^m}$ . In addition, we cannot access expected k-of-N losses directly; we must estimate them by sampling from  $\Psi$ . Therefore, the selection of a single  $\tau$  is a "best arm identification" problem, where  $\mathcal{T}^m$  is the set of arms. The Successive Rejects (SR) [Audibert et al., 2010] algorithm is an exploration-only bandit algorithm that can be used to solve this problem with a worst-case guarantee on the probability that it identifies the best arm. The more SR iterations that are run, the more likely it is to select the best arm. Algorithm 1 shows how to implement RPOSST for the sequential move model using regret matching<sup>+</sup> for tuning the test case weights and SR for the final selection of an m-tuple.

In specific applications, an example of which we will see in Section 4.2 and our experiments, we can construct our optimization game so that it is deterministic, and consequently, we can replace SR with a simple argmax.

The  $RPOSST_{\text{SEQ}}$  objective is the percentile performance loss

$$\min_{\substack{\tau \in \mathcal{T}^m \\ \hat{\sigma}_{\tau} \in \Delta^m}} \inf_{\substack{y \in \mathcal{Y} \\ \eta \in [0,1] \\ \mathbb{P}[\ell(\hat{\sigma}_{\tau}; \pi_{j}, \sigma) \leq y(\eta)] \geq \eta}} \int_{\eta \in [0,1]} y(\eta) \mu_{k\text{-of-}N}(d\eta), \tag{2}$$

where  $\langle \pi_i, \sigma \rangle \sim \Psi$ .

The sequential move model represents the policy testing use case where we select and fix m test cases and test case weights for all future test policies. Test scores are easily reproducible and comparable across test applications since the test cases never change.

**Theorem 4.1.** After  $T' \sim \text{Unif}(\{1,\ldots,T_1\})$ ,  $T_1 > 0$ , rounds of its optimization game, Algorithm I selects an m-tuple of test cases,  $\tau^*$  and weights  $\hat{\sigma}_{\tau^*}^{T'} \in \triangle^m$  that, with probability  $(1-p)(1-q)(1-\alpha)$ ,  $p,q,\alpha>0$ , are  $\frac{\varepsilon}{q}$ -optimal for Equation (2), where  $\varepsilon = \mathcal{O}\left(\sqrt{\frac{1}{T_1}m} + \sqrt{\frac{1}{T_1}\log(1/p)}\right)$  and  $\alpha = \mathcal{O}(\mathrm{e}^{-T_2})$ .

All proofs deferred to the Appendix. In the extreme case where  $\Pi_{\text{TNP}}$  covers  $\Pi$ , then this optimality result, (in terms of an upper bounded percentile loss integral), extends to all deployment candidates  $\Pi_{\text{CDP}}$ .

#### 4.2 DETERMINISTIC CVAR RPOSST

While in general, an RPOSST algorithm has a randomized procedure and a non-deterministic optimality guarantee,

<sup>&</sup>lt;sup>1</sup>RPOSST is run for T' rather than  $T_1$  rounds because we cannot guarantee a decrease in worst-case loss after every round. See the proof of Theorem D.2 for more details.

we can actually select hyperparameters so that RPOSST is deterministic, making the procedure simpler and more reliable. If we fix the ratio  $^k/N$  and allow  $N \to \infty$ , the k-of-N robustness measure converges toward the CVaR measure at the  $^k/N$  fractile. A k-of-N algorithm where  $N \to \infty$  cannot be implemented with the usual sampling procedure, but it can be implemented if the distribution characterizing our uncertainty,  $\Psi$ , has finite support.

Sampling  $\Psi$  infinitely would result in sampling all tuning-policy–target-distribution pairs in its support exactly in proportion to their probabilities. Rather than selecting k tuning-policy–target-distribution pairs, the antagonist must select pairs until their cumulative probability sums to k/N. Effectively, the antagonist assigns weights

$$\alpha_{(i)} = \min \left\{ \Psi \left( \left\langle \pi_{j_{(i)}}, \sigma_{(i)} \right\rangle \right), k/N - \sum_{h=1}^{i-1} \alpha_{(h)} \right\}$$

to each tuning-policy–target-distribution pair in  $\Psi$ 's support, where the ordering between pairs is determined by the loss each induces for the protagonist. Finally, these tuning-policy–target-distribution pairs are sampled according to the normalized weights  $\frac{\alpha_{(i)}N}{k}$ .

The robustness guarantees become deterministic because the entire RPOSST algorithm, denoted as  $\text{CVaR}(\eta)$  RPOSST for the  $\eta = {}^k/N$  fractile, can be run using exact expectations (excluding randomness in A, which is taken as given in RPOSST). Determinism in RPOSST<sub>SEQ</sub> allows us to directly check the exact expected loss of each test case distribution on each round, letting us track the lowest loss test case distribution across all rounds. This tracking, in turn, allows us to avoid both sampling T' and running the SR algorithm to do the final selection. Instead, we can simply return the lowest loss test case distribution across all T rounds.

If there are d tuning-policy-target-distribution pairs in  $\Psi$ 's support, then the expected  $\mathrm{CVaR}(\eta)$  loss of the protagonist on round t is  $L^t = \min_{\tau \in \mathcal{T}^m} \sum_{i=1}^d \frac{\alpha_{(i)}}{\eta} \ell(\hat{\sigma}_{\tau}^t; \pi_{j_{(i)}, \sigma_{(i)}})$ . The round with the lowest expected loss is  $t^* = \arg\min_{t \in \{1, \dots, T\}} L^t$ , and this definition allows us to state the following corollary.

**Corollary 4.2.** Assume that  $\Psi \in \triangle^d$  for some finite  $d \geq 1$ . After T rounds of the  $CVaR(\eta)$  RPOSST<sub>SEQ</sub> optimization game, where the protagonist chooses m-size tests according to regret matching<sup>+</sup> against a best response antagonist,  $\tau^*$  and  $\sigma_{\tau^*}^{t^*}$  are  $\varepsilon$ -optimal for Equation (2) under the  $\eta$ -fractile CVaR robustness measure, where  $\varepsilon = \mathcal{O}\left(\sqrt{\frac{1}{T}m}\right)$ .

Pseudocode for  $\text{CVaR}(\eta)$  RPOSST<sub>SEQ</sub> is presented in Appendix Algorithm 2.

In addition, we can construct a series of ablations of CVaR RPOSST $_{\text{SEQ}}$  to act as baselines for experiments, and to make a connection to the test-construction literature.

CVaR RPOSST<sub>SEQ</sub> generalizes an intuitive algorithm: find the m-tuple of test cases that minimizes the maximum error assuming a uniform distribution over the tuple. This minimax uniform algorithm is implemented by executing only the initialization and selection steps of CVaR(0) RPOSST<sub>SEQ</sub> (T=0). Further simplifying, minimax(TTD) uniform performs the antagonist maximization only over target distributions and assumes a uniform distribution over tuning policies. Minimax(TNP) uniform performs the antagonist maximization only over tuning policies and assumes a uniform target distribution. Miniaverage uniform assumes both a uniform distribution over tuning policies and for the target distribution.

Additionally, we could select test cases one at a time to minimize the maximum error, echoing greedy algorithms from the test-construction literature (Chapter 4 of van der Linden [2005]). This *iterative minimax* algorithm is almost the same as running the initialization and return steps of CVaR(0) RPOSST<sub>SEQ</sub> to select a single test case in a loop. The sole difference being that iterative minimax could select the same test case multiple times within its loop to adjust the test case weighting away from uniform.

#### 5 EXPERIMENTS

We explore CVaR RPOSST<sub>SEQ</sub>'s performance in three twoplayer games spanning the range of complexity from a toy one-shot game to a high-fidelity racing simulator, in comparison with minimax and miniaverage baselines. We show that robustness does tend to decrease test score errors on holdout policies and that RPOSST specifically either outperforms or performs about as well as each baseline in each domain.

#### 5.1 EXPERIMENTAL SETUP

In each domain, we start with data from playing out every pairing of n>0 policies, yielding a matrix of scores for the column policy. Each policy along the rows of this matrix is then treated as a test case, making the score at row i and column j the result of evaluating policy j on test case i.

To emulate unknown deployment candidate policies to be tested, we hold out h>0 columns of this matrix and call the policy associated with a holdout column a holdout policy. The remaining columns represent the test case results for the set of tuning policies. The resulting  $n\times (n-h)$  matrix is shifted and rescaled so that all entries are between zero and one, and then it is set as the test result matrix A that our methods take as input. Note, although h test cases are generated by holdout policies, as test cases they cannot provide any special information about what tests would be effective on the holdout policies. To simulate scenarios where the set of tuning policies covers the set of future candidate deployment policies to varying degrees, we run

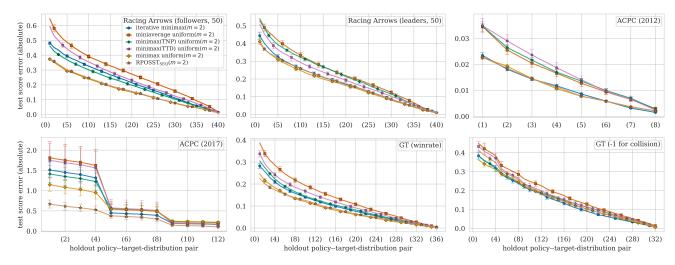


Figure 2: Expected test score error (absolute difference) across holdout-policy-target-distribution pairs on (top left and middle) Racing Arrows, (top right) the 2012 two-player, limit competition of the ACPC, (bottom left) the 2017 two-player, no-limit competition of the ACPC, (bottom middle and right) Gran Turismo<sup>TM</sup> 7 races, between CVaR(1%) RPOSST<sub>SEQ</sub> and baseline tests on 100 randomly sampled sets of holdout policies (20% of the full set of policies; 80% used as tuning policies). Holdout-policy-target-distribution pairs are sorted according to test score error. Each RPOSST<sub>SEQ</sub> instance was run for 500 rounds (T=500). Errorbars represent 95% t-distribution confidence intervals.

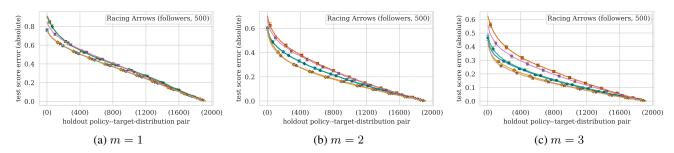


Figure 3: Expected test score error across holdout-policy–target-distribution pairs on Racing Arrows where test cases are follower policies. Here, 500 Racing Arrows policies were sampled for both the follower and leader role and then 96% of policies of both roles were held out before running RPOSST and each baseline. Each column uses a different test size m. 100 sets of holdout policies were sampled and each RPOSST<sub>SEO</sub> instance was run for 500 rounds (T = 500).

experiments with three different values of h: 0.2n, 0.4n, and 0.6n. 100 different holdout sets are randomly sampled for each value of h and in each domain.

Given results for n test cases, the goal is to produce a distribution over m < n test cases that provides accurate test results on the set of holdout policies, according to a set of target distributions. For our experiments, we use  $m \in \{1,2,3\}$  and the set of target distributions generated from the softmax function applied to the negative average test case result under four different scales, specifically,  $\exp\left(\frac{-\beta}{n}A\mathbf{1}\right)/\mathbf{1}^{\top}\exp\left(\frac{-\beta}{n}A\mathbf{1}\right)$  for  $\beta \in \{0,1,2,4\}$ , so that the distributions put varying degrees of emphasis on test cases that are more difficult on average across the tuning policies. We set the RPOSST uncertainty distribution,  $\Psi$ , to be uniform over each tuning-policy-target-distribution pair. We set the CVaR percentile to 1% so that it is nearly optimizing for the worst-case, but is slightly less pessimistic, to add

an additional distinguishing factor to RPOSST compared to the minimax and minaverage baselines. We use the absolute difference loss for both optimization and evaluation.

#### 5.2 DOMAINS

We test RPOSST on the following three domains of varying complexity. Each domain has two variants arising from asymmetry, multiple datasets, or alternative scoring rules. Appendix Appendix G provides further details on each domain.

**Racing Arrows.** Racing Arrows is a two-player, zero-sum, one-shot, continuous action game invented for our experiments to replicate aspects of a passing scenario in a race featuring a "leader" player and faster "follower" player. The follower tries to pass the leader while the latter tries to block. Scores are recorded as 0 or +1 for a loss or win, respec-

tively, for the column player, which is either the leader or the follower, depending on the configuration. We run RPOSST on both configurations. For our experiments, we sample 50 or 500 different leader and follower policies evenly spread through the valid policy space, angles in  $[0,\pi]$ , by taking 50 or 500 evenly spaced angles between  $[0.05\pi, (1-0.05)\pi]$  and then shifting them independently with uniform samples in  $[-0.05\pi, 0.05\pi]$ .

Annual Computer Poker Competition. We take two open datasets from the Annual Computer Poker Competition (ACPC) [Bard et al., 2013] containing pairwise match data for poker agents submitted to the 2017 two-player, no-limit competition and the 2012 two-player, limit competition. These competitions contain different agent populations since they are separated by five years and are in different game formats (limit and no-limit). The 2017 competition consists of 15 agents and the 2012 competition consists of 12 agents. Scores are recorded as chip differentials of duplicate matches (two sets of hands where players play with the same set of shuffled decks in both seats).

**Gran Turismo**<sup>TM</sup> **one-on-one races.** Gran Turismo<sup>TM</sup> 7 (GT)<sup>2</sup> is a high fidelity racing simulator on the PlayStation<sup>TM</sup> platform. Previous versions of GT served as benchmarks for training RL policies [Fuchs et al., 2021, Song et al., 2021] including policies that outraced the best human competitors [Wurman et al., 2022] in four-on-four racing. We consider a simpler one-on-one racing scenario (see Appendix Appendix G.3 for details). We carry out two experiments, one where test case results are average winrates, and another where policies receive 0 for a loss, +1 for a win, and -1 if there was a collision, making the game non-zero-sum. The test case pool is comprised of 43 trained RL policies and 3 built-in "AI" policies.

### 5.3 RESULTS AND ANALYSIS

The results of running CVaR(1%) RPOSST<sub>SEQ</sub> on each domain, with m=2 and 20% of policies marked as holdout policies, are shown in Figure 2. The same set of figures with m=1 and m=3, as well as 40% and 60% holdout policies, are qualitatively similar, except that the differences between the algorithms are typically smaller, and are provided in Appendix Appendix G.4.

Looking across each domain and variant, we can see that RPOSST<sub>SEQ</sub> performs nearly as well or better than all of the minimax and miniaverage baselines, particularly in terms of maximum error across holdout-policy–target-distribution pairs. Interestingly, RPOSST<sub>SEQ</sub> has noticeably lower error in ACPC 2017 and GT (winrate) on the four most difficult holdout-policy–target-distribution pairs to accurately evaluate. The improvement over the next best method is substantial in ACPC 2017 because RPOSST is the only method

with an unlimited ability to optimize with a non-uniform test case weighting. On the other variant in each domain, RPOSST $_{SEQ}$  is within the group of the lowest error methods. In the two Racing Arrows domains, RPOSST $_{SEQ}$  and minimax uniform substantially outperform the other methods, at least on the most difficult holdout-policy–target distribution pairs. This result shows that robustness is indeed beneficial here, but the uniform distribution over the selected two opponents happens to be quite effective. The GT variant where -1 is assigned to a collision appears to be more difficult than the winrate variant, as all the methods cluster together in this variant at higher errors than in the winrate variant.

These results illustrate the utility of incorporating robustness generally, as all of the robust methods tend to outperform miniaverage uniform. Minimax uniform and iterative minimax are the only baselines that minimize their maximum error over both tuning policy and target distribution uncertainty, and they are usually the next best methods after RPOSST<sub>SEQ</sub>. Minimax(TNP) uniform typically outperforms minimax(TTD) uniform, showing that it is more important to be robust to the tuning policy than the target distribution, in these domains. When the target distributions are the same in the optimization and holdout evaluation phases, robustness should directly improve the minimum performance across holdout realizations. Since no effort was made to enforce any relationship between the tuning and holdout policies, this result suggests that robustness to the tuning policy can yield large error reductions when  $\Pi_{TNP}$  are even somewhat similar to the holdout policies.

As an example of RPOSST's capabilities, consider the pairs of opponent policies chosen as test cases in GT (winrate) over 100 experiment seeds (Appendix Table 2). RPOSST<sub>SEO</sub> is both more accurate (Figure 2) and very consistent, choosing the same pair 90% of the time. Figure 4 illustrates the portion of the result matrix for just the two test cases most frequently chosen by RPOSST<sub>SEQ</sub> (test races against opponents 16 and 41). The race against policy 41 (bottom row) is chosen because that policy wins/loses about half the time, providing a 50/50 information split. Policy 16 is a weaker policy in many ways (more blue in the top row) but it serves to differentiate the worst policies (darker red squares in the left side of the matrix) from the rest of the policies, and to highlight the strongest policies. Specifically, the best performing policies almost always win against policy 16, which provides a strong complementary signal to the noisier but more competitive policy 41 test case. Overall, the two test cases indicate policies 1, 29, and 43 (darkest blue columns) are the strongest for deployment. Policy 1 is a built-in AI in an overpowered car but 29 and 43 are very strong RL policies. Looking at the overall winrate matrix (Appendix Figure 5b) we see that the same conclusion (the three dark-

<sup>2</sup>https://www.gran-turismo.com/us/

<sup>&</sup>lt;sup>3</sup>Iterative minimax can change its test case distribution away from uniform, but only indirectly by selecting a test case it already selected on a previous iteration before it fills its test-case quota.

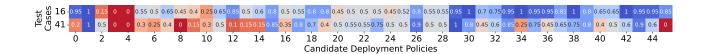


Figure 4: The GT test results of candidate deployment policies against the test case pair most favoured by RPOSST. Blue and red indicates positive and negative winrates respectively for the candidate deployment policy.

est blue columns overall) would have been chosen using all 46 test cases. Compressing from 46 test cases to two presents a massive saving in test time for future policies, and shows RPOSST<sub>SEQ</sub> can construct small tests to select deployment policies in a real and complex video game.

The results in Figure 3 repeat the previous analysis in Racing Arrows but with ten times the number of policies. Only the results where follower policies are treated as test cases are shown, but the corresponding results where leader policies are test cases appear similar and are shown in Appendix Appendix G.4. 96% of policies are held out, including those used as test cases, so there are only 20 test cases and tuning policies for RPOSST and the other algorithms to utilize. This experiment emulates a scenario where an efficient test is constructed once with a relatively small number of tuning policies and then reused for many future deployment candidates. As in the previous experiments, RPOSST is almost always one of the best methods.

#### 6 RELATED WORK

The bulk of the work on policy selection in RL focuses on selecting opponents for *training* with self-play algorithms [Hernandez et al., 2021]. In that case, diversity is key for training additional policies to cooperate [Rahman et al., 2022] or compete [Liu et al., 2021, McAleer et al., 2022] with pre-existing policies. However, the selection of policies as training opponents is often guided by aggregate performance metrics across entire populations [Li et al., 2019, Lanctot et al., 2017, Omidshafiei et al., 2019, Balduzzi et al., 2018] and thus do not reduce the number of opponent pairings (test cases) required for assessments.

On the testing side, researchers in complex domains develop procedures for testing skill competency using hand-calibrated [Wurman et al., 2022] or randomly generated tests with complex percentile-scoring functions [Team et al., 2021]. Our work seeks to automate and target test construction in such scenarios. Complementary work [Rowland et al., 2019] treats the computation of a result matrix as a multi-armed bandit problem, each entry represented by one arm. While this method can greatly reduce sampling costs in the presence of low-variance outcomes, it does not generalize to policies outside its input population, with the testing of a new policy requiring adding extra arms to be estimated from

scratch. However, this method could be used in tandem with RPOSST to reduce the samples required to compute A.

Learning to rank methods [Oosterhuis and de Rijke, 2021, Bruch, 2021, Hu et al., 2018] aim to find a function that ranks a set of items (*e.g.*, documents) based on the relevance of a given query, with hopes to generalize to future queries. Indeed, Akiyama et al. [2016] use learning to rank to evaluate action sequences. However, predicting unseen policy performances under this model requires the tuning policies to be the queries, which would produce a ranking of the test cases themselves. The scores from such tests would therefore be incomparable across policies, violating one of our main objectives.

Test construction in educational modeling [van der Linden, 2005] starts from an item bank and a statistical model (e.g., Item Response Theory [Embretson et al., 2000]) predicting the probability of answering each item correctly given a student's (unobserved) skill level. That model yields an *information matrix* and then automatic test construction methods, including linear optimization or greedy heuristics, can then build a finite-sized test. By contrast, we do not assume a model of the response variance or a univariate skill measurement, so a closed-form calculation of information is often infeasible. However, we do empirically compare our optimization approach to the greedy heuristic.

#### 7 CONCLUSION AND FUTURE WORK

RPOSST is, to the best of our knowledge, the first algorithm to directly address test construction for reinforcement learning policies. By leveraging the k-of-N framework, RPOSST provides bounds on the approximation error of the resulting test despite uncertainty over the exact policies that will be evaluated and the desired test case weighting in the future. Thus, RPOSST provides a much needed tool for policy selection in real-world deployment scenarios. An interesting direction for future work is generating the test cases themselves [Marris et al., 2021, Pugh et al., 2016], which is challenging on its own [Balduzzi et al., 2019].

#### Acknowledgements

Thanks to Francesco Riccio for reviewing this work. Thanks to the whole Sony AI team for experiment infrastructure.

#### References

- Hidehisa Akiyama, Masashi Tsuji, and Shigeto Aramaki. Learning evaluation function for decision making of soccer agents using learning to rank. In 2016 Joint 8th International Conference on Soft Computing and Intelligent Systems (SCIS) and 17th International Symposium on Advanced Intelligent Systems (ISIS), pages 239–242. IEEE, 2016.
- Jean-Yves Audibert, Sébastien Bubeck, and Rémi Munos. Best arm identification in multi-armed bandits. In *COLT*, pages 41–53. Citeseer, 2010.
- David Balduzzi, Karl Tuyls, Julien Perolat, and Thore Graepel. Re-evaluating evaluation. *Advances in Neural Information Processing Systems*, 31, 2018.
- David Balduzzi, Marta Garnelo, Yoram Bachrach, Wojciech Czarnecki, Julien Perolat, Max Jaderberg, and Thore Graepel. Open-ended learning in symmetric zero-sum games. In *International Conference on Machine Learn*ing, pages 434–443. PMLR, 2019.
- Nolan Bard, John Hawkin, Jonathan Rubin, and Martin Zinkevich. The annual computer poker competition. *AI Magazine*, 34(2):112–112, 2013.
- Noam Brown and Tuomas Sandholm. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science*, 359(6374):418–424, 2018.
- Noam Brown and Tuomas Sandholm. Superhuman AI for multiplayer poker. *Science*, 365(6456):885–890, 2019.
- Sebastian Bruch. An alternative cross entropy loss for learning-to-rank. In *Proceedings of the Web Conference* 2021, pages 118–126, 2021.
- Neil Burch. *Time and space: Why imperfect information games are hard.* PhD thesis, University of Alberta, 2017.
- Abraham Charnes and William W Cooper. Chance-constrained programming. *Management science*, 6(1): 73–79, 1959.
- Katherine Chen and Michael Bowling. Tractable objectives for robust policy optimization. *Advances in Neural Information Processing Systems*, 25:2069–2077, 2012.
- S.E. Embretson, S.E. Embretson, and S.P. Reise. *Item Response Theory for Psychologists*. Multivariate applications book series. L. Erlbaum Associates, 2000. ISBN 9780805828184.

- Meta FAIR, Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. Human-level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378(6624): 1067–1074, 2022.
- Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Stochastic regret minimization in extensive-form games. In *International Conference on Machine Learning*, pages 3018–3028, 2020.
- Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55 (1):119–139, 1997.
- Florian Fuchs, Yunlong Song, Elia Kaufmann, Davide Scaramuzza, and Peter Dürr. Super-human performance in Gran Turismo Sport using deep reinforcement learning. *IEEE Robotics and Automation Letters*, 6(3):4257–4264, 2021. doi: 10.1109/LRA.2021.3064284.
- Daniel Hernandez, Kevin Denamganai, Sam Devlin, Spyridon Samothrakis, and James Alfred Walker. A comparison of self-play algorithms under a generalized framework. *IEEE Transactions on Games*, 14(2):221–231, 2021.
- Yujing Hu, Qing Da, Anxiang Zeng, Yang Yu, and Yinghui Xu. Reinforcement learning to rank in e-commerce search engine: Formalization, analysis, and application. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 368–377, 2018.
- Michael Johanson, Nolan Bard, Neil Burch, and Michael Bowling. Finding optimal abstract strategies in extensive form games. In *26th AAAI Conference on Artificial Intelligence (AAAI-12)*, 2012.
- Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent reinforcement learning. *Advances in neural information processing systems*, 30, 2017.
- Ang Li, Ola Spyra, Sagi Perel, Valentin Dalibard, Max Jaderberg, Chenjie Gu, David Budden, Tim Harley, and Pramod Gupta. A generalized framework for population based training. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1791–1799, 2019.
- Xiangyu Liu, Hangtian Jia, Ying Wen, Yujing Hu, Yingfeng Chen, Changjie Fan, Zhipeng Hu, and Yaodong Yang. Towards unifying behavioral and response diversity for open-ended learning in zero-sum games. *Advances in Neural Information Processing Systems*, 34:941–952, 2021.

- Edward Lockhart, Marc Lanctot, Julien Pérolat, Jean-Baptiste Lespiau, Dustin Morrill, Finbarr Timbers, and Karl Tuyls. Computing approximate equilibria in sequential adversarial games by exploitability descent. In *IJCAI* 2019, 2019a.
- Edward Lockhart, Marc Lanctot, Julien Pérolat, Jean-Baptiste Lespiau, Dustin Morrill, Finbarr Timbers, and Karl Tuyls. Computing approximate equilibria in sequential adversarial games by exploitability descent. *arXiv* preprint arXiv:1903.05614, 2019b.
- Luke Marris, Paul Muller, Marc Lanctot, Karl Tuyls, and Thore Graepel. Multi-agent training beyond zero-sum with correlated equilibrium meta-solvers. In *International Conference on Machine Learning*, pages 7480–7491. PMLR, 2021.
- Stephen McAleer, Kevin Wang, John B Lanier, Marc Lanctot, Pierre Baldi, Tuomas Sandholm, and Roy Fox. Anytime psro for two-player zero-sum games. 2022.
- Colin McDiarmid. Concentration. In *Probabilistic methods* for algorithmic discrete mathematics, pages 195–248. 1998.
- Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisỳ, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.
- Shayegan Omidshafiei, Christos Papadimitriou, Georgios Piliouras, Karl Tuyls, Mark Rowland, Jean-Baptiste Lespiau, Wojciech M Czarnecki, Marc Lanctot, Julien Perolat, and Remi Munos. α-rank: Multi-agent evaluation by evolution. *Scientific reports*, 9(1):1–29, 2019.
- Harrie Oosterhuis and Maarten de de Rijke. Robust generalization and safe query-specialization counterfactual learning to rank. In *Proceedings of the Web Conference* 2021, pages 158–170, 2021.
- Julien Perolat, Bart De Vylder, Daniel Hennes, Eugene Tarassov, Florian Strub, Vincent de Boer, Paul Muller, Jerome T Connor, Neil Burch, Thomas Anthony, et al. Mastering the game of stratego with model-free multiagent reinforcement learning. *Science*, 378(6623):990– 996, 2022.
- Justin K Pugh, Lisa B Soros, and Kenneth O Stanley. Quality diversity: A new frontier for evolutionary computation. *Frontiers in Robotics and AI*, page 40, 2016.
- Arrasy Rahman, Elliot Fosong, Ignacio Carlucho, and Stefano V Albrecht. Towards robust ad hoc teamwork agents by creating diverse training teammates. *arXiv* preprint *arXiv*:2207.14138, 2022.

- Mark Rowland, Shayegan Omidshafiei, Karl Tuyls, Julien Perolat, Michal Valko, Georgios Piliouras, and Remi Munos. Multiagent evaluation under incomplete information. *NeurIPS*, 32, 2019.
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez,
  Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc
  Lanctot, Sander Dieleman, Dominik Grewe, John Nham,
  Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap,
  Madeleine Leach, Koray Kavukcuoglu, Thore Graepel,
  and Demis Hassabis. Mastering the game of go with
  deep neural networks and tree search. *Nature*, 529(7587):
  484–489, 2016.
- Yunlong Song, HaoChih Lin, Elia Kaufmann, Peter Dürr, and Davide Scaramuzza. Autonomous Overtaking in Gran Turismo Sport Using Curriculum Reinforcement Learning. In *ICRA*, 2021.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Oskari Tammelin. Solving large imperfect information games using cfr+. *arXiv preprint arXiv:1407.5042*, 2014.
- Oskari Tammelin, Neil Burch, Michael Johanson, and Michael Bowling. Solving heads-up limit texas hold'em. In 24th International Joint Conference on Artificial Intelligence (IJCAI 2015), 2015.
- Open Ended Learning Team, Adam Stooke, Anuj Mahajan, Catarina Barros, Charlie Deck, Jakob Bauer, Jakub Sygnowski, Maja Trebacz, Max Jaderberg, Michaël Mathieu, Nat McAleese, Nathalie Bradley-Schmieg, Nathaniel Wong, Nicolas Porcel, Roberta Raileanu, Steph Hughes-Fitt, Valentin Dalibard, and Wojciech Marian Czarnecki. Open-ended learning leads to generally capable agents. *ArXiv Pre-Print*, abs/2107.12808, 2021. URL https://arxiv.org/abs/2107.12808.
- Wim J. van der Linden. *Linear models for optimal test design*. Springer, 2005.
- Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, pages 1–5, 2019.
- Peter R Wurman, Samuel Barrett, Kenta Kawamoto, James MacGlashan, Kaushik Subramanian, Thomas J Walsh, Roberto Capobianco, Alisa Devlic, Franziska Eckert, Florian Fuchs, et al. Outracing champion gran turismo drivers with deep reinforcement learning. *Nature*, 602 (7896):223–228, 2022.

### A APPENDIX

## **B** GLOSSARY

**Policy.** A policy to solve a control problem or play a game, potentially generated by an RL algorithm.

**Deployment policy.** A policy used in production, *e.g.*, deployed to end users, used in a competition, or integrated into a technology demonstration.

**Deployment candidate.** A policy in consideration for deployment.

**Test.** The aggregate result of test cases applied to a policy.

**Test case.** An atomic unit of a test that reveals a particular skill or emulates a specific deployment scenario. RPOSST selects a small number of test cases and a distribution over them so that we can avoid executing all conceivable test cases on every deployment candidate every time we want to deploy a policy.

**Test case result.** The numerical result of evaluating a policy on a test case. This number should be a good estimate of the policy's expected performance in the test case scenario, but it maybe noisy if the test case is stochastic, *e.g.*, the average test case result observed from Monte Carlo rollouts.

**Test score.** The final score produced by a test, *i.e.*, the average test case result across test cases, perhaps weighted by the relative importance of each test case.

**Tuning policy.** A policy used at the start of the RPOSST procedure to gather information about test cases. Each tuning policy is evaluated on each test case to construct the test case result matrix that forms the basis of RPOSST's loss function.

#### C THEORY BACKGROUND

We make use of six basic results, which are restated here for completeness.

**Proposition C.1** (Azuma-Hoeffding inequality.). For constants  $\langle c^t \in \mathbb{R} \rangle_{t=1}^T$ , martingale difference sequence  $\langle Y^t \in \mathbb{R} \rangle_{t=1}^T$  where  $|Y^t| \leq c^t$  for each t, and  $\tau \geq 0$ ,

$$\mathbb{P}\left[\left|\sum_{t=1}^T Y^t\right| \geq \tau\right] \leq 2 \exp\left(\frac{-\tau^2}{2\sum_{t=1}^T (c^t)^2}\right).$$

For proof, see that of Theorem 3.14 by McDiarmid [1998].

**Proposition C.2** (Regret matching<sup>+</sup> regret bound). Consider an online decision process with m actions and the set of bounded, linear loss functions,  $\mathcal{L} = [0, L]^m$ . Regret matching<sup>+</sup> accumulates pseudoregrets  $q^{1:t} = [q^{1:t-1} + \rho^t]_+$ ,  $q^{1:0} = \mathbf{0}$ , where  $\rho^t = (\ell^t)^\top \sigma^t - \ell^t$  is the instantaneous regret on round t under loss function  $\ell^t \in \mathcal{L}$ , and  $\sigma^t = q^{1:t-1}/(\mathbf{1}^\top q^{1:t-1})$  if  $\mathbf{1}^\top q^{1:t-1} > 0$  or  $\sigma^t = \frac{1}{m}\mathbf{1}$  otherwise, is regret matching<sup>+</sup>'s action distribution on round t. After T rounds, regret matching<sup>+</sup>'s cumulative regret is bounded as  $\sum_{t=1}^T \rho^t \leq L\sqrt{Tm}$ .

For proof, see Tammelin et al. [2015].

**Proposition C.3** (The linearization trick). Consider an online decision process with convex decision set  $\Theta \subseteq \mathbb{R}^m$  and a set of bounded, convex loss functions  $\mathcal{L} \subseteq \{\ell \mid \ell : \Theta \to [0, L]\}$ , where each loss function  $\ell \in \mathcal{L}$  has subgradients with bounded maximum magnitude, i.e.,  $\|\nabla \ell(\theta)\|_{\infty} \leq G$ , for all  $\theta \in \Theta$ . The instantaneous regret under loss function  $\ell \in \mathcal{L}$  is upper bounded by the instantaneous regret under the loss function subgradient  $\nabla \ell(\theta)$  given decision  $\theta \in \Theta$ , i.e.,

$$\ell(\theta) - \ell(\theta') \le (\nabla \ell(\theta))^{\top} \theta - (\nabla \ell(\theta))^{\top} \theta'.$$

*Proof.* From the convexity of  $\ell$ , its first-order Taylor expansion lower bounds  $\ell$ , *i.e.*,  $\ell(\theta') \ge \ell(\theta) + (\nabla \ell(\theta))^{\top}(\theta' - \theta)$ , for all  $\theta, \theta' \in \Theta$ . Therefore,

$$\ell(\theta) - \ell(\theta') \le \ell(\theta) - \left(\ell(\theta) + (\nabla \ell(\theta))^{\top} (\theta' - \theta)\right)$$
$$= (\nabla \ell(\theta))^{\top} \theta - (\nabla \ell(\theta))^{\top} \theta',$$

as required.

**Proposition C.4** (Lemma 2 of Lockhart et al. [2019a,b]). Assume that on each round t of an online decision process with decision set  $\Theta \subseteq \mathbb{R}^m$  and bounded loss functions from  $\mathcal{L} \subseteq \{\ell \mid \ell : \Theta \to [0, L]\}$ , the loss function  $\ell^t$  maximizes the loss of  $\theta^t \in \Theta$  chosen by the decision-maker, i.e.,  $\ell^t \in \arg\max_{\ell \in \mathcal{L}} \ell(\theta^t)$ . On the round  $t^*$  where the minimum loss was observed,  $t^* \in \arg\min_{t \in \{1, \dots, T\}} \ell^t(\theta^t)$ , the decision  $\theta^{t^*}$  has a maximum loss that is no more than  $\frac{1}{T}\rho^{1:T}(\theta)$  larger than that of any alternative decision  $\theta \in \Theta$ , i.e.,  $\ell^{t^*}(\theta^{t^*}) - \ell^{\theta}(\theta) \leq \frac{1}{T}\rho^{1:T}(\theta)$ , where  $\ell^{\theta}(\theta) \in \mathcal{L}$  is a loss function that maximizes the loss on  $\theta$ .

*Proof.* Since the loss function on each round is chosen to maximize loss, the average regret for not choosing  $\theta \in \Theta$  is lower bounded as

$$\frac{1}{T}\rho^{1:T} \ge \frac{1}{T} \min_{t \in \{1, ..., T\}} T\ell^{t}(\theta^{t}) - \frac{1}{T} \sum_{t=1}^{T} \ell^{t}(\theta)$$

$$\ge \ell^{t^{*}}(\theta^{t^{*}}) - \ell^{\theta}(\theta),$$

as required.

**Proposition C.5** (Theorem 4 of Johanson et al. [2012]). Assume that on each round t of an online decision process with decision set  $\Theta \subseteq \mathbb{R}^m$  and bounded (possibly random) loss functions from  $\mathcal{L} \subseteq \{\ell \mid \ell : \Theta \to [0, L]\}$ , the loss function  $\ell^t$  maximizes the loss of  $\theta^t \in \Theta$  chosen by the decision-maker, i.e.,  $\ell^t \in \arg\max_{\ell \in \mathcal{L}} \ell(\theta^t)$ . The loss function that the decision-maker observes on each round t may be a random loss function  $\hat{\ell}^t$  where  $\mathbb{E}\left[\hat{\ell}^t\right] = \ell^t$ . On round  $T' \sim \mathrm{Unif}(\{1, \ldots, T\})$  after T rounds of the online decision process, the decision  $\theta^{T'}$  has a maximum loss that is no more than  $\frac{1}{qT}\rho^{1:T}(\theta)$  larger than that of any alternative decision  $\theta \in \Theta$  with probability 1-q,  $q \in (0,1]$ , i.e.,  $\ell^{T'}(\theta^{T'}) - \ell^{\theta}(\theta) \leq \frac{1}{qT}\rho^{1:T}(\theta)$  holds with probability 1-q, where  $\ell^{\theta}(\theta) \in \mathcal{L}$  is a loss function that maximizes the loss on  $\theta$  and the cumulative regret  $\rho^{1:T}$  is with respect to the expected loss functions,  $\langle \ell^t \rangle_{t=1}^T$ .

See Johanson et al. [2012] for proof.

**Proposition C.6** (Successive Rejects error probability). Consider a best action identification task with m actions from set A. Each time an action  $a \in A$  is selected, a random sample of that action's loss,  $\ell(a) \in [-0.5, 0.5]$ , under a fixed but random loss function  $\ell$ , is observed. The goal is to identify an action  $a^* \in A^* \subset A$  with the lowest expected loss,  $\mathbb{E}[\ell(a^*)]$ , after T samples. The probability that the action returned by the Successive Rejects algorithm is in  $A^*$  is at least

$$1 - \frac{m(m-1)}{2} \exp\left(-\frac{T-m}{\overline{\log}(m)H_2}\right),\,$$

where  $\overline{\log}(m) = \frac{1}{2} + \sum_{i=2}^{m} \frac{1}{i}$ ,  $H_2 = \max_{i \in \{1, \dots, |\mathcal{A} \setminus \mathcal{A}^*|\}} \frac{i}{\left(\mathbb{E}\left[\ell(a_{(i)})\right] - \mathbb{E}\left[\ell(a^*)\right]\right)^2}$ , and  $a_{(i)}$  is the action that achieves the  $i^{th}$  smallest loss (with ties broken arbitrarily) among the suboptimal actions.

See Audibert et al. [2010] for proof.

## D SEQUENTIAL-MOVE MODEL THEORY

**Lemma D.1.** Consider a k-of-N game with m actions and the set of bounded, convex loss functions  $\mathcal{L} = \{\ell \mid \ell : \triangle^m \to [0,L]\}$ , where each loss function  $\ell \in \mathcal{L}$  has subgradients with bounded maximum magnitude, i.e.,  $\|\nabla \ell(\sigma)\|_{\infty} \leq G$ , for all  $\sigma \in \triangle^m$ . Let the k-worst loss functions from N of those sampled from the given uncertainty distribution  $\Psi$  on round t be  $\langle \ell_{(i)}^t \in \mathcal{L} \rangle_{i=1}^k$ . The randomly sampled k-of-N loss function on round t is then the average  $\bar{\ell}^t = \frac{1}{k} \sum_{i=1}^k \ell_{(i)}$ . After T rounds, regret matching<sup>+</sup> on the random loss gradients  $\nabla \bar{\ell}^t(\sigma^t)$  has no more than  $2G\sqrt{Tm} + 2L\sqrt{2T\log^{1}/p}$  cumulative regret on the expected k-of-N losses,  $\langle \mathbb{E}[\bar{\ell}]^t \rangle_{t=1}^T$ , with probability 1-p, p>0.

*Proof.* Since regret matching<sup>+</sup> observes and learns directly from  $\nabla \bar{\ell}^t$ , its regret for not always choosing  $\sigma \in \triangle^m$ , under the sampled loss functions, is deterministically upper bounded as

$$R^{1:T} = \sum_{t=1}^{T} \underline{\bar{\ell}^t(\sigma^t) - \bar{\ell}^t(\sigma)} \le 2G\sqrt{Tm},$$

where  $\langle \sigma^t \in \triangle^m \rangle_{t=1}^T$  are the decisions made by regret matching<sup>+</sup>. This bound comes from regret matching<sup>+</sup>'s regret bound on linear losses (Proposition C.2) and the linearization trick (Proposition C.3), which states that the regret on loss gradients upper bounds that of the loss itself, i.e.,  $R^{1:T} \leq \sum_{t=1}^T \left( \nabla \bar{\ell}^t(\sigma^t) \right)^\top \sigma^t - \left( \nabla \bar{\ell}^t(\sigma) \right)^\top \sigma$ .

The rest of the proof largely follows the proof of Farina et al. [2020]'s Proposition 1. The sequence of differences,  $\langle \mathbb{E}[R^t] - R^t \leq 2L \rangle_{t=1}^T$ , is a bounded martingale difference sequence.

The probability that the expected cumulative regret,  $\mathbb{E}[R^{1:T}]$ , is bounded by the cumulative sampled regret plus slack  $\lambda \geq 0$  is bounded according to the Azuma-Hoeffding inequality (Proposition C.1) as

$$\mathbb{P}\left[\mathbb{E}[R^{1:T}] \le R^{1:T} + \lambda\right] \tag{3}$$

$$\leq \mathbb{P}\left[\sum_{t=1}^{T} \mathbb{E}[R^t] - R^t \leq \lambda\right] \tag{4}$$

$$=1-\mathbb{P}\left[\sum_{t=1}^{T}\mathbb{E}[R^{t}]-R^{t}\geq\lambda\right] \tag{5}$$

$$\leq 1 - \exp\left(\frac{2\lambda^2}{4T(2L)^2}\right).$$
(6)

Setting  $\lambda = 2L\sqrt{2T\log(1/p)}$  ensures that

$$\mathbb{E}[R^{1:T}] \le R^{1:T} + 2L\sqrt{2T\log 1/p}$$

with probability 1 - p. Since  $R^{1:T} \leq 2L\sqrt{Tm}$ ,

$$\mathbb{E}[R^{1:T}] \le 2G\sqrt{Tm} + 2L\sqrt{2T\log 1/p}$$

with probability 1 - p, as required.

**Theorem D.2.** After  $T' \sim \mathrm{Unif}(\{1,\ldots,T_1\})$ ,  $T_1 > 0$ , rounds of its optimization game, Algorithm 1 selects an m-tuple of test cases,  $\tau^*$  and weights  $\hat{\sigma}_{\tau^*}^{T'} \in \triangle^m$  that, with probability  $(1-p)(1-q)(1-\alpha)$ ,  $p,q,\alpha>0$ , are  $\frac{\varepsilon}{q}$ -optimal for Equation (2), where  $\varepsilon = \mathcal{O}\Big(\sqrt{\frac{1}{T_1}m} + \sqrt{\frac{1}{T_1}\log(1/p)}\Big)$  and  $\alpha = \mathcal{O}\big(\mathrm{e}^{-T_2}\big)$ .

*Proof.* Recall that the k-of-N loss  $\bar{\ell}^t$  that RPOSST<sub>SEQ</sub> updates from on each round  $t=1,\ldots,T_1$  is a Monte Carlo estimate of the k-of-N percentile loss,

$$L_{\mu_{k\text{-of-}N},\Psi}(\hat{\sigma}_{\tau}^{t}) = \inf_{y \in \mathcal{Y}} \int_{\substack{\eta \in [0,1] \\ \mathbb{P}_{\pi_{j},\sigma}\left[\ell(\hat{\sigma}_{\tau}^{t};\pi_{j},\sigma) \leq y(\eta)\right] \geq \eta}} y(\eta)\mu_{k\text{-of-}N}(d\eta) = \mathbb{E}_{\pi_{j},\sigma}\left[\bar{\ell}^{t}\right], \tag{7}$$

where  $(\pi_j, \sigma) \sim \Psi$ . The sequence of test case weights,  $\langle \sigma_{\tau}^t \rangle_{t=1}^{T_1}$ , for each m-tuple of test cases  $\tau \in \mathcal{T}$  is therefore random. All of the following probabilities and expectations are with respect to these random variables.

Lemma D.1 guarantees that RPOSST<sub>SEQ</sub>, in generating the test case weight sequence  $\langle \sigma_{\tau}^t \rangle_{t=1}^{T_1}$  has no more than  $C = 2G\sqrt{T_1m} + 2L\sqrt{2T_1\log 1/p}$  cumulative regret on the k-of-N percentile losses,

$$\rho_{\sigma_{\tau}}^{1:T_{1}} = \sum_{t=1}^{T} L_{\mu_{k \text{-of-}N}, \Psi}(\hat{\sigma}_{\tau}^{t}) - L_{\mu_{k \text{-of-}N}, \Psi}(\sigma_{\tau}),$$

for not always selecting test case weights  $\sigma_{\tau}$ , with probability 1-p. That is,  $1-p=\mathbb{P}\left[\rho_{\sigma_{\tau}}^{1:T_1}\leq C\right]$ .

Proposition C.5 guarantees that, on round  $T' \sim \mathrm{Unif}(1,\ldots,T_1)$ , the weights for each m-tuple are  $\frac{1}{qT_1}\rho_{\sigma_{\tau}}^{1:T_1}$  close to optimal for Equation (7), with probability 1-q. That is,  $1-q=\mathbb{P}\bigg[L_{\mu_{k\text{-of-}N},\Psi}(\hat{\sigma}_{\tau}^{T'})-L_{\mu_{k\text{-of-}N},\Psi}(\sigma_{\tau})\leq \frac{\rho_{\sigma_{\tau}}^{1:T_1}}{qT_1}\bigg]$ ,

$$\text{and this holds regardless of the value of} \quad \rho_{\sigma_{\tau}}^{1:T_{1}}, \quad i.e., \quad \mathbb{P}\left[L_{\mu_{k\text{-of-}N},\Psi}(\hat{\sigma}_{\tau}^{T'}) - L_{\mu_{k\text{-of-}N},\Psi}(\sigma_{\tau}) \leq \frac{\rho_{\sigma_{\tau}}^{1:T_{1}}}{qT_{1}}\right] \quad = \\ \mathbb{P}\left[L_{\mu_{k\text{-of-}N},\Psi}(\hat{\sigma}_{\tau}^{T'}) - L_{\mu_{k\text{-of-}N},\Psi}(\sigma_{\tau}) \leq \frac{\rho_{\sigma_{\tau}}^{1:T_{1}}}{qT_{1}} \mid \rho_{\sigma_{\tau}}^{1:T_{1}} \leq C'\right] \text{ for all } C' \in \mathbb{R}.$$

Combining these two results, we see that the probability that  $\hat{\sigma}_{\tau}^{T'}$  has at most  $\frac{C}{qT_1}$  excess k-of-N percentile loss is

$$\mathbb{P}\left[L_{\mu_{k \text{-of-}N}, \Psi}(\hat{\sigma}_{\tau}^{T'}) - L_{\mu_{k \text{-of-}N}, \Psi}(\sigma_{\tau}) \leq \frac{C}{qT_{1}}\right] = \mathbb{P}\left[L_{\mu_{k \text{-of-}N}, \Psi}(\hat{\sigma}_{\tau}^{T'}) - L_{\mu_{k \text{-of-}N}, \Psi}(\sigma_{\tau}) \leq \frac{\rho_{\sigma_{\tau}}^{1:T_{1}}}{qT_{1}}, \rho_{\sigma_{\tau}}^{1:T_{1}} \leq C\right] \\
= \mathbb{P}\left[L_{\mu_{k \text{-of-}N}, \Psi}(\hat{\sigma}_{\tau}^{T'}) - L_{\mu_{k \text{-of-}N}, \Psi}(\sigma_{\tau}) \leq \frac{\rho_{\sigma_{\tau}}^{1:T_{1}}}{qT_{1}} \mid \rho_{\sigma_{\tau}}^{1:T_{1}} \leq C\right] \mathbb{P}\left[\rho_{\sigma_{\tau}}^{1:T_{1}} \leq C\right] \\
(9)$$

$$= \mathbb{P}\left[L_{\mu_{k\text{-of-}N},\Psi}(\hat{\sigma}_{\tau}^{T'}) - L_{\mu_{k\text{-of-}N},\Psi}(\sigma_{\tau}) \le \frac{\rho_{\sigma_{\tau}}^{1:T_{1}}}{qT_{1}}\right] \mathbb{P}\left[\rho_{\sigma_{\tau}}^{1:T_{1}} \le C\right]$$

$$\tag{10}$$

$$= (1-p)(1-q). (11)$$

The last remaining step is to complete the outer minimization in Equation (2) to select a single m-tuple of test cases. Since the k-of-N loss observed on each round is random, we cannot compute a simple argmin using the test case weights on round T', and are instead faced with a best arm identification problem. For this, we run the Successive Rejects algorithm, which we know from Proposition C.6 identifies a minimum loss m-tuple of test cases with probability at least

$$\alpha = 1 - \frac{m(m-1)}{2} \exp\left(-\frac{T_2 - m}{\overline{\log}(m)H_2}\right).$$

The probability of selecting the best m-tuple using the test case weights on round T' is independent of whether or not the regret bound C was actually achieved or if the test case weights on T' are actually nearly optimal for any given m-tuple, the probability of which we previously characterized as (1-p)(1-q). Therefore, the probability of achieving  $\frac{C}{qT_1}$ -optimality given each m-tuple and selecting the best m-tuple is the product  $(1-p)(1-q)(1-\alpha)$ , as required.

The  $\sqrt{|\mathcal{T}|}$  dependence in Theorem D.2 could be improved to  $\sqrt{\log(|\mathcal{T}|)}$  if regret matching<sup>+</sup> (within or without CFR, respectively) was replaced with an algorithm like Hedge [Freund and Schapire, 1997], but this tends to lead to worse performance in practice (see, *e.g.*, Tammelin et al. [2015], Burch [2017]).

In the deterministic CVaR RPOSST case, we get the following corollary.

**Corollary D.3.** Assume that  $\Psi \in \triangle^d$  for some finite  $d \geq 1$ . After T rounds of the  $CVaR(\eta)$  RPOSST<sub>SEQ</sub> optimization game, where the protagonist chooses m-size tests according to regret matching<sup>+</sup> against a best response antagonist,  $\tau^*$  and  $\sigma_{\tau^*}^{t^*}$  are  $\varepsilon$ -optimal for Equation (2) under the  $\eta$ -fractile CVaR robustness measure, where  $\varepsilon = \mathcal{O}\left(\sqrt{\frac{1}{T}m}\right)$ .

*Proof.* Proposition C.2 and Proposition C.4 ensures that there is a round  $t_{\tau}^* \leq T$  where  $\sigma_{\tau}^{t_{\tau}^*}$  is  $2G\sqrt{m\frac{1}{T}}$ -optimal on the deterministic k-of-N losses. Since the k-of-N loss function observed on each round is deterministic, we can perform a simple minimization across  $\{1,\ldots,T\}$  and the m-tuple of test cases to find the minimizers  $t^*$  and  $\tau^*$ , leading to the stated optimality guarantee.  $\Box$ 

## E DETERMINISTIC CVAR $(\eta)$ RPOSST<sub>SEO</sub> PSEUDOCODE

Pseudocode for  $CVaR(\eta)$  RPOSST<sub>SEQ</sub> is presented in Algorithm 2.

### F SIMULTANEOUS-MOVE MODEL

We present a more in-depth description of the simultaneous move antagonist model which describes the RPOSST $_{SIM}$  as introduced in Section 4. This description is complemented by pseudocode describing its workings in Algorithm 3.

In this model, the antagonist does not observe which m-tuple of test cases,  $\tau$ , is sampled from the protagonist's  $\hat{\sigma}_{\mathcal{T}}^t \in \triangle^{|\mathcal{T}|^m}$  distribution, making the antagonist role more difficult. The simultaneous move model corresponds to the policy testing use case where a new m-tuple of test cases is sampled independently for each test that is performed. Effectively, the protagonist and antagonist choose  $\tau$  and  $\left\langle \left\langle \pi_{j_{(i)}}, \sigma_{(i)} \right\rangle \right\rangle_{i=1}^k$  respectively in a simultaneous fashion. In this model, the antagonist must choose a single list of tuples  $\left\langle \left\langle \pi_{j_{(i)}}, \sigma_{(i)} \right\rangle \right\rangle_{i=1}^k$  that will lead to a large loss across all of the m-tuples of test cases that the protagonist might choose, thereby preventing the antagonist from exploiting the lacking aspects of each individual m-tuple.

The protagonist in the simultaneous move model must carefully choose  $\hat{\sigma}_{\mathcal{T}}^t$  and each m-tuple distribution,  $[\hat{\sigma}_{\tau}^t]_{\tau \in \mathcal{T}^m}$ , to thwart the antagonist. We organize the protagonist's actions into two sequential decisions: first choosing the m-tuple  $\tau$  and then choosing  $\hat{\sigma}_{\tau}^t$  given  $\tau$ . We then use CFR<sup>+</sup> to refine both  $\hat{\sigma}_{\tau}^t$  and each  $\hat{\sigma}_{\tau}^t$  after each round.

Instantiating the percentile performance loss of Equation (1) for the simultaneous move model, the RPOSST objective is,

$$\min_{\substack{\hat{\sigma}_{\mathcal{T}} \in \triangle^{|\mathcal{T}|^m} \\ [\hat{\sigma}_{\tau} \in \triangle^m]_{\tau \in \mathcal{T}^m}}} \inf_{y \in \mathcal{Y}} \int_{\eta \in [0,1]} y(\eta) \mu_{k\text{-of-}N}(d\eta), \tag{12}$$

where  $\sigma, \pi_i \sim \Psi$ .

After (linearly) averaging the protagonist's choices of  $\hat{\sigma}_{\mathcal{T}}^t$  and  $[\hat{\sigma}_{\tau}^t]_{\tau \in \mathcal{T}^m}$  across each round, Algorithm 3 returns the average distributions  $\bar{\sigma}_{\mathcal{T}}^T$  and  $[\bar{\sigma}_{\tau}^T]_{\tau \in \mathcal{T}^m}$ .

The simultaneous-move model can be made deterministic using a CVaR measure in the same way as the sequential-move model. If we fix the ratio k/N and allow  $N \to \infty$ , the k-of-N robustness measure converges toward the CVaR measure at the k/N fractile. Furthermore, if our the distribution characterizing our uncertainty,  $\Psi$ , is over a discrete set of manageable size, then we can run RPOSST on CVaR robustness measures. In RPOSST<sub>SIM</sub>, the lowest loss test case distributions across all rounds can also be tracked instead of averaging all of the distributions.

## **G** EXPERIMENTAL DETAILS

In this section we provide further details on some of the experimental setups used in Section 5.

All CVaR(1%) RPOSST<sub>SEQ</sub> procedures were run on a 16 core AMD<sup>®</sup> Ryzen 7 5800h CPU with 30.7 GiB of memory. See Table 1 for the time required to run CVaR(1%) RPOSST<sub>SEQ</sub> on each domain.

#### G.1 RACING ARROWS

Racing Arrows is a two-player, zero-sum, one-shot, continuous action game that replicates simple aspects of a passing scenario in a race featuring a "leader" player and faster "follower" player. The goal of the follower is to pass the leader while the goal of the leader is to block the follower.

Both players privately choose an angle in the half-circle between 0 and pi for their arrow. The speed of each player is represented as the length of their arrow. The leader and follower are assigned a speed according to their roles, where the leader's speed of 0.8 is slightly slower than the follower's speed of 1 to give the follower a chance to pass. The distance a player travels is the height of their arrow, *i.e.*, speed  $\cdot \sin(\text{angle})$ .

The follower is blocked and the leader wins if the difference between the two arrows is below  $\pi/10$ , that is, the leader is close enough to block the follower. If the follower is not blocked, then the player who traveled the farthest wins. Players receive +1 for a win, 0 for a loss, or 0.5 if they travel exactly the same distance (these payoffs sum to the constant +1, which is isomorphic to true zero-sum payoffs).

#### G.2 ANNUAL COMPUTER POKER COMPETITION

The Annual Computer Poker Competition (ACPC) was run to test autonomous poker playing agents from 2006 to 2017. The logs of play are freely available online.<sup>4</sup> Typically, these competitions are Texas hold'em variants: two-player limit,

<sup>&</sup>lt;sup>4</sup>http://www.computerpokercompetition.org/downloads/competitions

two-player no-limit, and 3-player limit, where "limit" and "no-limit" indicates whether players are only allowed to bet in fixed increments or if they can bet any number of chips from their current stack, respectively. Chip stacks reset to their initial sizes after every hand (Doyle's game) so that players can be evaluated on their average one-hand performance across deck shufflings and seat positions.

To reduce variance, hands are played in duplicate, which means that the same deck order is played out multiple times so that each player has a turn playing with the same hands. For example, if Alice in seat 1 is dealt the ace and king of spades and Bob in seat 2 is dealt the 2 and 7 of hearts in one hand, then Alice and Bob will also play the same hand in opposite positions, where Bob is dealt the ace and king of spades in seat 1, and Alice is dealt the 2 and 7 of hearts. Alice's duplicate score is then the number of chips she wins over what Bob won in the same position, averaged across both positions.

Our experiments use duplicate score data, *i.e.*, a test case result here is a duplicate score between two agents, from the 2012 two-player limit and the 2017 two-player no-limit events.

## G.3 GRAN TURISMO<sup>TM</sup> 7

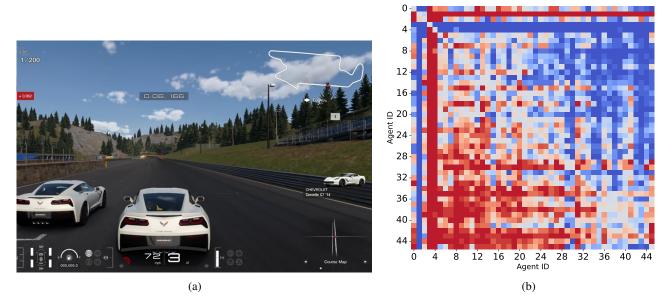


Figure 5: Figure 5a shows a screenshot of two RL agents racing at the Trial Mountain racetrack. The layout of the track can be seen in the top right. Figure 5b shows the result matrix for the zero-sum experiment. Blue / red colors indicate positive / negative winrates from the point of view of the column player. Agents 0-2 correspond to the different built-in AIs, with the remaining agents being the trained RL agents sorted according to skill. Diagonal values denote an agent playing against itself, which we artificially set to 50%.

Our Gran Turismo<sup>TM</sup> 7 experiments were conducted using the Gran Turismo<sup>TM</sup> 7 racing simulator. Previous versions of the Gran Turismo<sup>TM</sup> 7 franchise have been used to exhibit reinforcement learning results [Fuchs et al., 2021, Song et al., 2021] including outracing top human drivers [Wurman et al., 2022]. Note our focus was not on agent training but rather the problem of selecting the best policy for a deployment, so for training we used the same training parameters reported by by Wurman et al. except for changes to training scenarios to match the track and car combination chosen for this experiment, training only for one-on-one competition, and utilizing a version of self-play to simplify the training process.

The experiment was conducted at the Trial Mountain racetrack (see Figure 5a) with the RL policy (and any RL-trained opponent policies) driving a Chevrolet Corvette C7 Stingray '14 using Sport Hard tires. The track and car were chosen because the long straightaways and sharp turns at Trial Mountain led to competitive racing among various RL policies as

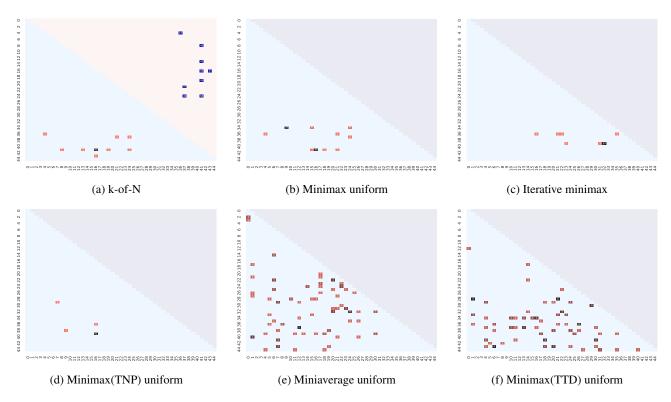


Figure 6: Triangular matrices denoting frequencies of test case pairs choosen by all 6 algorithmic ablations over 100 runs of the winrate GT experiment on a holdout of size 20%. This visualization is possible because only 2 test cases where chosen as output. The upper triangular matrix from Figure 6a denotes average probability mass given to test case i. All other algorithms are limited to uniformly mixing over test cases so the upper triangular matrix is omitted for clarity.

there are many different areas of the track where passes can occur and the long straightaways allow the agent to use the slipstream of the other car to stay in touch with the car in front.

From a single one-on-one training run we evaluated checkpoints from epochs 5, 200, and then every 75 epochs between epoch 1000 and 4000 for a total of 43 checkpoints. We also evaluated 3 built-in AI agent using cars and tires that made them competitive with the RL agents. Overall we evaluated 46 policies, each of which was considered as a candidate deployment policy or an opponent in a test case.

To create the result matrix shown in Figure 5b, each race was run 20 times with a side-by-side standstill start with the candidate and opponent policies swapping sides half the time to enforce symmetry. An agent would obtain 1 or 0 for winning or losing the race respectively. The diagonal denoting a race between an agent against itself was filled in with 0.5 entries. As a second experiment on Gran Turismo<sup>TM</sup> 7 for a non-zero sum game, using the sportsmanship rule mentioned in Section 5 we recomputed the result matrix from Figure 5b so as to penalize trajectories where any car collisions had happened, giving both agents a payoff of -1. We remove the entries in the result matrix related to built-in AIs as they are highly collision averse and therefore the sportsmanship constraints would not change their test results, reducing the test case pool size to 43.

### G.4 SUPPLEMENTAL EXPERIMENTAL RESULTS

Figure 2 in Section 5 analyses the quantitative performance of RPOSST and its algorithmic ablations with respect to measuring test scores on a holdout set of unseen candidate deployment policies. We complement that analysis with a qualitative study of behaviors exhibited by the algorithms using the large GT experiment with holdout of size 20 as a representative example. We are interested in examining (1) how deterministic each algorithm's output is with respect to the selection of test case pairs and (2) whether different algorithms choose the same test-cases.

The lower triangular matrices from Figure 6 show the frequency at which test case pairs were chosen over the 100 seeds. The top 2 most selected test case pairs for each algorithm are presented in Table 2. We observe that RPOSST, alongside Iterative minimax and Minimax(TNP) uniform are very deterministic algorithms, favouring the selection of the same test

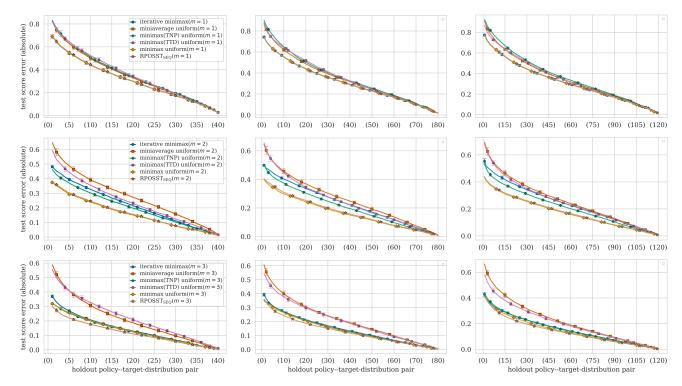


Figure 7: Expected test score error (absolute difference) across holdout-policy–target-distribution pairs on Racing Arrows where test cases are 50 follower policies. Each row uses a different setting for the test size (m=1 top, m=2 middle, and m=3 bottom) and each column uses a different holdout proportion (20% held out in the left column, 40% middle, and 60% right). 100 sets of holdout policies were sampled. Holdout-policy–target-distribution pairs are sorted according to test score error. Each RPOSST<sub>SEQ</sub> instance was run for 500 rounds (T=500). Errorbars represent 95% t-distribution confidence intervals.

case pair over 90%, 87% and 92% of the seeds respectively. We deem this a desirable property, as variance in evaluation scenarios is undesirable because it can hamper interpretability and reproducibility. In contrast, Minimax uniform exhibits a bimodal choice. The remaining algorithms feature a very high variance in their choice of test case pairs, with their most chosen test case pair being selected 5% of the time, spreading selection widely.

From Table 2, test case 16 is heavily favoured by half of the algorithms (RPOSST<sub>SEQ</sub>, Minimax uniform and Minimax(TNP) uniform), followed to a lesser extent by test case 41. This indicates that all these algorithms find useful structure in such pairs of agents.

In Figures 7 to 12, we show the performance of RPOSST<sub>SEQ</sub> and baselines in each domain across test sizes ( $m \in \{1, 2, 3\}$ ) and holdout proportions (20%, 40%, and 60%). Figure 13 shows the results for the 500 policy Racing Arrows experiment where the leader policies are treated as test cases.

We note that as m increases, the error on the holdout set typically decreases, particularly for RPOSST, since larger tests have the capacity to be strictly more accurate. A qualitative analysis of these results suggests that there are few substantial differences between RPOSST tests of different sizes or with different, reasonably sized, holdout sets. Furthermore, the performance ordering of the tested algorithms remains the same as the results presented in the main paper.

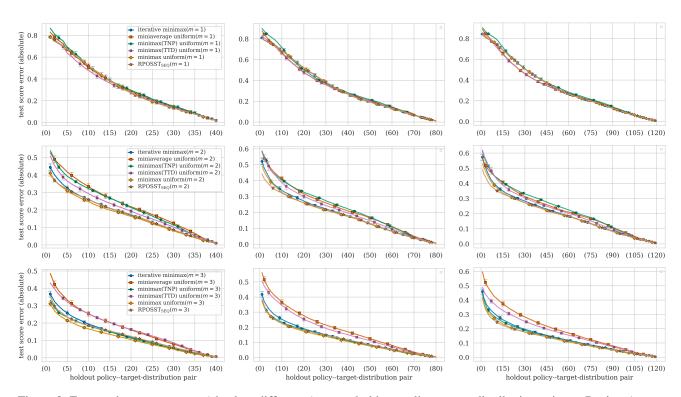


Figure 8: Expected test score error (absolute difference) across holdout-policy-target-distribution pairs on Racing Arrows where test cases are 50 leader policies. Each row uses a different setting for the test size (m=1 top, m=2 middle, and m=3 bottom) and each column uses a different holdout proportion (20% held out in the left column, 40% middle, and 60% right). 100 sets of holdout policies were sampled. Holdout-policy-target-distribution pairs are sorted according to test score error. Each RPOSST<sub>SEQ</sub> instance was run for 500 rounds (T=500). Errorbars represent 95% t-distribution confidence intervals.

```
1 Inputs: \langle \eta, T, m, \Psi, \tau^0, \ell \rangle
 2 q_{	au}^{1:0} \leftarrow \mathbf{0} \in \mathbb{R}^{m+|	au^0|} for 	au \in \mathcal{T}^m
 t^* \leftarrow 1
 4 \bar{v}^{t^*} \leftarrow -\infty
 5 for t \leftarrow 1, \dots, T do
            for \tau \in \mathcal{T}^m do
                    z^t \leftarrow \mathbf{1}^{\top} q_{\tau}^{1:t-1}
                    \hat{\sigma}_{	au}^t \leftarrow q_{	au}^{1:t-1}/z^t \text{ if } z^t > 0 \text{ else } 1/m
 8
                    // Fill in zeros so that \hat{\sigma}_{	au}^t \in \triangle^{|\mathcal{T}|}.
                    \hat{\sigma}_{\tau}^t(x) \leftarrow 0 \text{ for } x \in \mathcal{T} \setminus (\tau \cup \tau^0)
10
                    \left[\ell_{\tau,(i)}\right]_{i=1}^{d} \leftarrow \mathcal{L}_{\eta}(\hat{\sigma}_{\tau}^{t}, \eta, \Psi, \ell)
11
                   v_{\tau}^{t} \leftarrow -\sum_{i=1}^{d} \frac{\partial \ell_{\tau,(i)}, \pi_{j_{(i)}}}{\partial \hat{\sigma}_{\tau}^{t}}
                   // Update regret matching+.
                  \bar{v}_{\tau}^t \leftarrow (\hat{\sigma}_{\tau}^t)^{\top} v_{\tau}^t
                    \rho_{\tau}^{\tau} \leftarrow v_{\tau}^{t} - \bar{v}_{\tau}^{t} 
 q_{\tau}^{1:t} \leftarrow [q_{\tau}^{1:t-1} + \rho_{\tau}^{t}]_{+}
15
                     // Update the best round.
                    if \bar{v}_{	au}^t > \bar{v}^{t^*} then
18
                       t^* \leftarrow t
                         \bar{v}^{t^*} \leftarrow \bar{v}_{\tau}^t
21 return 	au^{t^*}, \hat{\sigma}_{	au^*}^{t^*}
 1 Procedure \mathcal{L}_{\eta} Inputs: \langle \hat{\sigma}, \eta, \Psi, \ell \rangle
              // The support of \Psi, \mathrm{supp}(\Psi), is assumed to be a finite number d=|\mathrm{supp}(\Psi)|.
             for \pi_{j_i}, \sigma_i \in \operatorname{supp}(\Psi) do
 3
                   // Evaluate \hat{\sigma}.
 4
                \ell_i \leftarrow \ell(\hat{\sigma}; \pi_{j_i}, \sigma_i)
 5
             \operatorname{Sort}\left(\left\{i\mid\ell_i\right\}_{i=1}^d\right)
 6
              // Assign weights to each loss function.
 7
             // Iterate over \Psi's support sorted accoding to descending loss value from the
                     previous step.
              \beta \leftarrow 0
 9
             for \pi_{j_{(i)}}, \sigma_{(i)} \in \operatorname{supp}(\Psi) do
10
                    \alpha_{(i)} = \min \left\{ \Psi \left( \left\langle \pi_{j_{(i)}}, \sigma_{(i)} \right\rangle \right), \eta - \beta \right\}
11
                 \beta \leftarrow \beta + \alpha_{(i)}
12
             return \left[\frac{\alpha_{(i)}}{\eta}\ell_{(i)}\right]_{i=1}^a
13
```

**Algorithm 2:** Deterministic  $CVaR(\eta)$  RPOSST<sub>SEO</sub> with regret matching<sup>+</sup>

```
1 Inputs: \langle k, N, T, m, \Psi, \tau^0, \ell \rangle
  2 // Initialize pseudoregrets.
 \mathbf{3} \ q_{\mathcal{T}}^{1:0} \leftarrow \mathbf{0} \in \mathbb{R}^{|\mathcal{T}^m|}
 4 q_{	au}^{1:0} \leftarrow \mathbf{0} \in \mathbb{R}^{m+|	au^0|} for 	au \in \mathcal{T}^m
  5 // Initialize average distributions.
  6 \hat{\sigma}_{\tau}^{1:0} \leftarrow \mathbf{0} \in \mathbb{R}^{|\mathcal{T}^m|}
 7 \hat{\sigma}_{\tau}^{1:0} \leftarrow \mathbf{0} \in \mathbb{R}^{m+|\tau^0|} for \tau \in \mathcal{T}^m
  s for t \leftarrow 1, \dots, T do
                    // Sample antagonist actions.
                    \pi_{i}, \sigma_{i} \sim \Psi for i = 1 \dots N
                    // Generate test case distributions.
11
                   z_{\mathcal{T}}^t \leftarrow \mathbf{1}^{\top} q_{\mathcal{T}}^{1:t-1}
12
                    \hat{\sigma}_{\mathcal{T}}^t \leftarrow q_{\mathcal{T}}^{1:t-1}/z_{\mathcal{T}}^t \text{ if } z_{\mathcal{T}}^t > 0 \text{ else } 1/|\mathcal{T}^m| \\ \text{for } \tau \in \mathcal{T}^m \text{ do} 
13
14
                            \begin{aligned} z_{\tau}^t &\leftarrow \mathbf{1}^{\top} q_{\tau}^{1:t-1} \\ \hat{\sigma}_{\tau}^t &\leftarrow q_{\tau}^{1:t-1} / z_{\tau}^t \text{ if } z_{\tau}^t > 0 \text{ else } \mathbf{1} / m \end{aligned} 
15
16
                        // Fill in zeros so that \hat{\sigma}_{\tau}^t \in \triangle^{|\mathcal{T}|}. \hat{\sigma}_{\tau}^t(x) \leftarrow 0 \text{ for } x \in \mathcal{T} \setminus (\tau \cup \tau^0)
17
18
                    // Evaluate the CFR<sup>+</sup> distributions.
19
                   \ell_i \leftarrow (\hat{\sigma}_{\mathcal{T}}^t)^{\top} [\ell(\hat{\sigma}_{\tau}^t; \pi_{j_i}, \sigma_i)]_{\tau \in \mathcal{T}^m} \text{ for } i = 1, \dots, N
20
                    // Sort to identify the worst k \, .
21
                    \operatorname{SortBy}\left(\left[\left\langle\sigma_{i},\pi_{j_{i}}\right\rangle\right]_{i=1}^{N},\left[\ell_{i}\right]_{i=1}^{N}\right)
22
                    // Update CFR+.
23
                    for 	au \in \mathcal{T}^m do
24
                           \begin{array}{l} \ell_{\tau,(i)} \leftarrow \ell(\hat{\sigma}_{\tau}^t; \pi_{j_{(i)}}, \sigma_{(i)}) \text{ for } i = 1, \dots, k \\ v_{\tau}^t \leftarrow \frac{-1}{k} \sum_{i=1}^k \frac{\partial \ell_{\tau,(i)}}{\partial \hat{\sigma}_{\tau}^t} \end{array}
25
26
                         \begin{aligned} \rho_{\tau}^t &\leftarrow v_{\tau}^t - (\hat{\sigma}_{\tau}^t)^{\top} v_{\tau}^t \\ q_{\tau}^{1:t} &\leftarrow [q_{\tau}^{1:t-1} + \rho_{\tau}^t]_{+} \end{aligned}
27
28
                  v_{\mathcal{T}}^{t} \leftarrow \frac{-1}{k} \sum_{i=1}^{k} \left[ \ell_{\tau,(i)} \right]_{\tau \in \mathcal{T}^{m}} \\ \rho_{\mathcal{T}}^{t} \leftarrow v_{\mathcal{T}}^{t} - (\hat{\sigma}_{\mathcal{T}}^{t})^{\mathsf{T}} v_{\mathcal{T}}^{t} \\ q_{\mathcal{T}}^{1:t} \leftarrow \left[ q_{\mathcal{T}}^{1:t-1} + \rho_{\mathcal{T}}^{t} \right]_{+}
30
31
                    // Update average distributions.
                  \hat{\sigma}_{\mathcal{T}}^{1:t} \leftarrow \hat{\sigma}_{\mathcal{T}}^{1:t-1} + t\hat{\sigma}_{\mathcal{T}}^{t} 
 \hat{\sigma}_{\tau}^{1:t} \leftarrow \hat{\sigma}_{\tau}^{1:t-1} + t\hat{\sigma}_{\mathcal{T}}^{t}(\tau)\hat{\sigma}_{\tau}^{t} \text{ for } \tau \in \mathcal{T}^{m}
```

**Algorithm 3:** RPOSST<sub>SIM</sub> (simultaneous model; CFR<sup>+</sup>)

Table 1: Approximate amount of time required to run T=500 rounds of CVaR(1%) RPOSST<sub>SEQ</sub> in each domain. Runtimes are similar across both variants in each domain and across holdout policy set sizes.

domain	runtime / seed	
Racing Arrows	$\sim 2$ minutes	
ACPC	$\sim 10$ seconds	
GT	$\sim 90$ seconds	

Table 2: Top two test case pairs and corresponding selection frequencies chosen by each algorithm over the 100 seeds in the large GT experiment.

Algorithm	Pairs	Frequency
$RPOSST_{SEQ}$	(41, 16)	90
	(41, 19)	3
Minimax uniform	(41, 16)	40
	(34, 9)	37
Iterative minimax	(39, 32)	87
	(39, 31)	3
Minimax(TNP) uniform	(40, 16)	92
	(36, 16)	4
Miniaverage uniform	(37, 12)	5
	(40, 1)	4
Minimax(TTD) uniform	(43, 6)	4
	(28, 1)	4

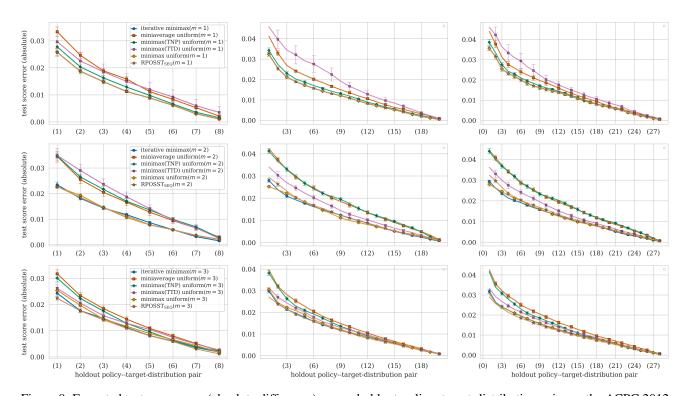


Figure 9: Expected test score error (absolute difference) across holdout-policy-target-distribution pairs on the ACPC 2012 data. Each row uses a different setting for the test size (m=1 top, m=2 middle, and m=3 bottom) and each column uses a different holdout proportion (20% held out in the left column, 40% middle, and 60% right). 100 sets of holdout policies were sampled. Holdout-policy-target-distribution pairs are sorted according to test score error. Each RPOSST<sub>SEQ</sub> instance was run for 500 rounds (T=500). Errorbars represent 95% t-distribution confidence intervals.

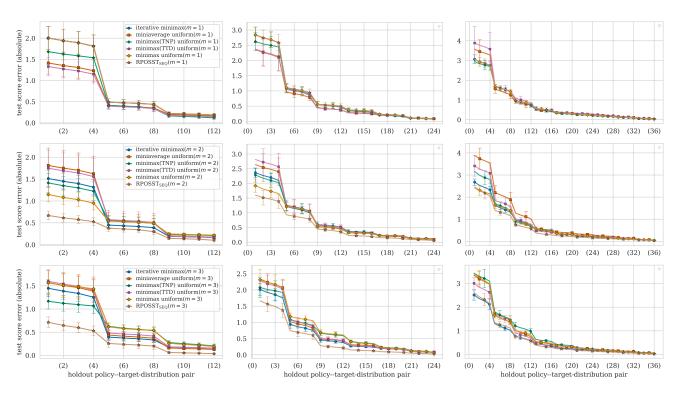


Figure 10: Expected test score error (absolute difference) across holdout-policy-target-distribution pairs on the ACPC 2017 data. Each row uses a different setting for the test size (m=1 top, m=2 middle, and m=3 bottom) and each column uses a different holdout proportion (20% held out in the left column, 40% middle, and 60% right). 100 sets of holdout policies were sampled. Holdout-policy-target-distribution pairs are sorted according to test score error. Each RPOSST<sub>SEQ</sub> instance was run for 500 rounds (T=500). Errorbars represent 95% t-distribution confidence intervals.

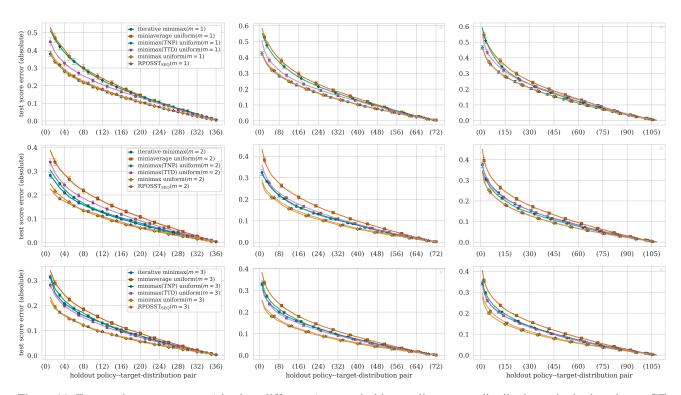


Figure 11: Expected test score error (absolute difference) across holdout-policy-target-distribution pairs in the winrate GT domain. Each row uses a different setting for the test size (m=1 top, m=2 middle, and m=3 bottom) and each column uses a different holdout proportion (20% held out in the left column, 40% middle, and 60% right). 100 sets of holdout policies were sampled. Holdout-policy-target-distribution pairs are sorted according to test score error. Each RPOSST<sub>SEQ</sub> instance was run for 500 rounds (T=500). Errorbars represent 95% t-distribution confidence intervals.

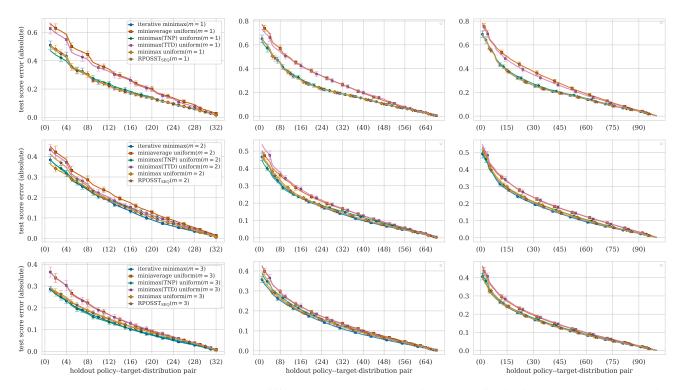


Figure 12: Expected test score error (absolute difference) across holdout-policy-target-distribution pairs in the GT domain where -1 is given for a collision. Each row uses a different setting for the test size (m=1 top, m=2 middle, and m=3 bottom) and each column uses a different holdout proportion (20% held out in the left column, 40% middle, and 60% right). 100 sets of holdout policies were sampled. Holdout-policy-target-distribution pairs are sorted according to test score error. Each RPOSST<sub>SEQ</sub> instance was run for 500 rounds (T=500). Errorbars represent 95% t-distribution confidence intervals.

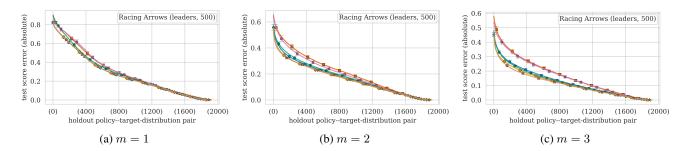


Figure 13: Expected test score error (absolute difference) across holdout-policy-target-distribution pairs on Racing Arrows where test cases are leader policies. Here, 500 Racing Arrows policies were sampled for both the follower and leader role and then 96% of policies of both roles were held out before running RPOSST and each baseline. Each column uses a different setting for the test size (m=1 top, m=2 middle, and m=3 bottom). 100 sets of holdout policies were sampled. Holdout-policy-target-distribution pairs are sorted according to test score error. Each RPOSST<sub>SEQ</sub> instance was run for 500 rounds (T=500). Errorbars represent 95% t-distribution confidence intervals.