

# Discovering Multiagent Learning Algorithms with Large Language Models

Zun Li  
Google DeepMind  
New York City, USA  
lizun@google.com

Daniel Hennes  
Google DeepMind  
Zurich, Switzerland  
hennes@google.com

John Schultz  
Google DeepMind  
New York City, USA  
jhtschultz@google.com

Marc Lanctot  
Google DeepMind  
Montreal, Canada  
lanctot@google.com

## ABSTRACT

Much of the advancement of Multi-Agent Reinforcement Learning (MARL) in imperfect-information games has historically depended on manual iterative refinement of baselines. While foundational families like Counterfactual Regret Minimization (CFR) and Policy Space Response Oracles (PSRO) rest on solid theoretical ground, the design of their most effective variants often relies on human intuition to navigate a vast algorithmic design space. In this work, we propose the use of AlphaEvolve, an evolutionary coding agent powered by large language models, to automatically discover new multiagent learning algorithms. We demonstrate the generality of this framework by evolving novel variants for two distinct paradigms of game-theoretic learning. First, in the domain of iterative regret minimization, we evolve the logic governing regret accumulation and policy derivation, discovering a new algorithm, Volatility-Adaptive Discounted (VAD-)CFR. VAD-CFR employs novel, non-intuitive mechanisms—including volatility-sensitive discounting, consistency-enforced optimism, and a hard warm-start policy accumulation schedule—to outperform state-of-the-art baselines like Discounted Predictive CFR+. Second, in the regime of population based training algorithms, we evolve training-time and evaluation-time meta strategy solvers for PSRO, discovering a new variant, Smoothed Hybrid Optimistic Regret (SHOR-)PSRO. SHOR-PSRO introduces a hybrid meta-solver that linearly blends Optimistic Regret Matching with a smoothed, temperature-controlled distribution over best pure strategies. By dynamically annealing this blending factor and diversity bonuses during training, the algorithm automates the transition from population diversity to rigorous equilibrium finding, yielding superior empirical convergence compared to standard static meta-solvers.

## KEYWORDS

Multi-Agent Reinforcement Learning, Game Theory, Large Language Models, Meta-Learning

### ACM Reference Format:

Zun Li, John Schultz, Daniel Hennes, and Marc Lanctot. 2026. Discovering Multiagent Learning Algorithms with Large Language Models. In *Proc. of the*

*Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026)*, C. Amato, L. Dennis, V. Mascardi, J. Thangarajah (eds.), May 25 – 29, 2026, Paphos, Cyprus. © 2026 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). This work is licensed under the Creative Commons Attribution 4.0 International (CC-BY 4.0) licence.

*25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026)*, Paphos, Cyprus, May 25 – 29, 2026, IFAAMAS, 27 pages.

## 1 INTRODUCTION

The field of Multi-Agent Reinforcement Learning has achieved remarkable milestones in recent years, reaching superhuman performance in domains ranging from Poker to real-time strategy games. These advances have been driven by a diverse array of methods, including game-theoretic regret minimization [4] and population-based league training [23]. The practical performance of these algorithms relies heavily on specific structural choices—such as how regret is discounted over time or how a specific equilibrium solution concept is derived. Historically, the refinement of these choices has been a largely manual endeavor. Researchers must rely on intuition and trial-and-error to navigate a vast combinatorial space of potential update rules, often defaulting to mathematically tractable heuristics (e.g., linear averaging or fixed discounting) that may not be optimal.

In this work, we propose to overcome this limitation by automating the design process itself with Large Language Models (LLMs). We apply AlphaEvolve [17], a distributed evolutionary system powered by LLMs, to the domain of multi-agent learning. Unlike traditional hyperparameter optimization or genetic programming, AlphaEvolve leverages the code-generation capabilities of LLMs to perform semantic evolution. By treating the algorithm’s source code as the genome, the system uses LLMs to act as intelligent genetic operators—performing mutation to rewrite logic, introduce new control flows, and inject novel symbolic operations. This allows the search to transcend simple parameter tuning and discover entirely new mechanisms for equilibrium finding.

We validate the generality of this framework by applying it to the two dominant paradigms of imperfect-information game solving:

1. Evolving Counterfactual Regret Minimization (CFR) [32]: The CFR family of algorithms relies on recursive definitions for accumulating regret and deriving an average policy. State-of-the-art variants like PCFR+ [8] and DCFR [3] were developed by manually identifying specific weighting schemes or regret targets to accelerate convergence. We apply AlphaEvolve to the symbolic operations governing regret accumulation, policy aggregation, and current policy derivation. The search yields Volatility-Adaptive Discounted (VAD-)CFR, a novel variant that dynamically adjusts its discounting parameters based on the volatility of regret updates and utilizes a

regret-magnitude weighted warm-start to construct the average strategy. This approach outperforms existing baselines across a diverse set of game benchmarks.

2. Evolving Policy-Space Response Oracles [12] (PSRO): PSRO generalizes the Double Oracle algorithm [13], expanding a population of policies by iteratively computing best responses and solving for a meta-strategy. Standard meta-solvers (e.g., Projected Replicator Dynamics or Uniform) enforce static trade-offs between exploration (expanding the game graph) and exploitation (refining the equilibrium). These static heuristics often fail to adapt to the changing topology of the empirical game during training. We apply AlphaEvolve to discover dynamic training-time and evaluation-time meta-solvers that optimize this schedule. The resulting variant, Smoothed Hybrid Optimistic Regret (SHOR-)PSRO, utilizes a hybrid update rule that mixes regret-based stability with aggressive greedy exploitation. By annealing the weight of these components over the training horizon, SHOR-PSRO empirically shows both algorithmic robustness and iteration efficiency.

The contributions of this paper are summarized as follows:

- (1) We showcase using LLM-driven evolution to design multi-agent learning algorithms, moving beyond parameter tuning to symbolic code evolution.
- (2) We identify VAD-CFR, an evolved algorithm that demonstrates how automated search can uncover update rules that are effective yet non-intuitive to human designers. In Appendix 7.2, we also present AOD-CFR, a variant discovered in an early trial that employs more conventional mechanisms while maintaining competitive performance.
- (3) We demonstrate the versatility of the framework by evolving PSRO meta-solvers that improve convergence speed and stability automating the transition from exploration to exploitation, yielding a new variant, SHOR-PSRO.

By automating the discovery of these mechanisms, we anticipate that future game-theoretic solvers will derive from a blend of human ingenuity and AI-driven insights.

## 2 GAME THEORETIC PRELIMINARIES

We formulate our problem within the framework of **Extensive-Form Games (EFGs)** with imperfect information, which models sequential interactions involving multiple agents and hidden information.

### 2.1 Extensive-Form Games and Exploitability

An  $N$ -player extensive-form game is  $\Gamma = \langle \mathcal{N}, \mathcal{H}, \mathcal{Z}, \mathcal{A}, u, \mathcal{I}, \sigma \rangle$ , where  $\mathcal{N} = \{1, \dots, N\}$  denotes the set of players [20].  $\mathcal{H}$  is the set of all possible histories (sequences of actions), where  $\mathcal{Z} \subseteq \mathcal{H}$  represents terminal histories. For  $i \in \mathcal{N}$ , let  $\mathcal{H}_i \subseteq \mathcal{H}$  represent the subset of histories where player  $i$  acts. For any non-terminal history  $h$ ,  $\mathcal{A}(h)$  is the set of legal actions. The utility function  $u_i : \mathcal{Z} \rightarrow \mathbb{R}$  assigns a payoff to player  $i$  at a terminal node.

Crucially, imperfect information is modeled via **Information Sets**  $I \in \mathcal{I}_i$ . Specifically,  $\mathcal{I}_i$  partitions the histories  $\mathcal{H}_i$  such that player  $i$  cannot distinguish between histories  $h, h' \in I$  (e.g., due to hidden cards). It is required that  $\mathcal{A}(h) = \mathcal{A}(h')$  for all  $h, h' \in I$ , so for simplicity we denote the legal actions at  $I$ ,  $\mathcal{A}(I)$ . A **strategy** (or policy)  $\sigma_i(I)$  assigns a probability distribution over actions  $a \in$

$\mathcal{A}(I)$  for each information set. A strategy profile  $\sigma = (\sigma_1, \dots, \sigma_N)$  determines the expected utility  $u_i(\sigma)$ . A **Nash Equilibrium (NE)** is a strategy profile  $\sigma^*$  such that neither player can increase their utility by deviating unilaterally:

$$u_i(\sigma_i^*, \sigma_{-i}^*) \geq u_i(\sigma'_i, \sigma_{-i}^*) \quad \forall \sigma'_i, \forall i \in \mathcal{N} \quad (1)$$

To measure the performance of our evolved algorithms, we use **Exploitability**. The exploitability of a strategy profile  $\sigma$  is the average of the incentives for players to deviate to their Best Response (BR):

$$\text{Exploitability}(\sigma) = \frac{1}{N} \sum_{i \in \mathcal{N}} \left( \max_{\sigma'_i} u_i(\sigma'_i, \sigma_{-i}) - u_i(\sigma) \right) \quad (2)$$

A strategy is an  $\epsilon$ -Nash Equilibrium if the exploitability is less than or equal to  $\epsilon$ . In small or medium-sized games, we can compute this value exactly by traversing the full game tree.

### 2.2 Counterfactual Regret Minimization (CFR)

CFR is an iterative algorithm that minimizes **counterfactual regret** [32]. It decomposes the global regret minimization problem into independent local regret minimization problems at each information set.

Let  $\pi^\sigma(h)$  be the probability of reaching history  $h$  under strategy  $\sigma$ . We define  $\pi_{-i}^\sigma(h)$  as the contribution of all players *except*  $i$  (including chance) to reaching  $h$ . The **counterfactual value** of reaching information set  $I$  and playing action  $a$  is the expected utility given that player  $i$  played to reach  $I$ :

$$v_i(\sigma, I, a) = \sum_{h \in I} \pi_{-i}^\sigma(h) \sum_{z \in \mathcal{Z}, h \subseteq z} \pi^\sigma(z | h, a) u_i(z) \quad (3)$$

The **instantaneous counterfactual regret** for not playing action  $a$  at iteration  $t$  is the difference between the counterfactual value of action  $a$  and the expected value at  $I$ :

$$r^t(I, a) = v_i(\sigma^t, I, a) - \sum_{a' \in \mathcal{A}(I)} \sigma^t(I, a') v_i(\sigma^t, I, a') \quad (4)$$

Standard CFR accumulates these values linearly over iterations  $T$ :

$$R^T(I, a) = \sum_{t=1}^T r^t(I, a) \quad (5)$$

The current policy  $\sigma^{t+1}$  is derived from accumulated regret, typically using **Regret Matching (RM)**, which assigns probabilities proportional to positive regret:

$$\sigma^{t+1}(I, a) = \frac{\max(R^t(I, a), 0)}{\sum_{a'} \max(R^t(I, a'), 0)} \quad (6)$$

The strategy  $\sigma^t$  at any single iteration may not converge to NE. Instead, CFR outputs the **average strategy**  $\bar{\sigma}^T$ , computed by weighting the iteration strategy  $\sigma^t$  by the player's contribution to the reach probability  $\pi_i^{\sigma^t}(I)$ .

**CFR Variants:** Several variants modify these update rules. For example, **CFR+** [22] replaces the regret accumulation  $R^T$  with floor bounding  $\max(R^{T-1} + r^t, 0)$  and uses linear averaging weights ( $w_t = t$ ) for the average policy. Our work uses AlphaEvolve to search for optimal variations of these accumulation and derivation functions.

### 2.3 Policy Space Response Oracles (PSRO)

PSRO [1, 12] acts as a meta-solver that generalizes the Double Oracle algorithm [13]. It operates on a higher level of abstraction called the **Meta-Game** (or Empirical Game [25]). PSRO maintains a population of policies  $\Pi_i = \{\sigma_i^1, \dots, \sigma_i^k\}$  for each player. The meta-game is represented by a payoff tensor  $M$ , where entries  $M_i^{j_1 \dots j_N} = u_i(\sigma_1^{j_1}, \dots, \sigma_N^{j_N})$  correspond to the expected utility of player  $i$  when policies from the population are pitted against each other.

**The PSRO Loop:** At each epoch  $k$ , the algorithm performs three steps:

- (1) **Meta-Strategy Solver (MSS) at training-time:** A solver computes a meta-strategy  $\phi_i$  (a probability distribution over the population  $\Pi_i$ ). Common solvers include **Uniform**  $\phi_i^j = 1/|\Pi_i|$  and **Nash**  $\phi$  which is a Nash Equilibrium of the current meta-game  $M$ .
- (2) **Oracle (Best Response):** A new policy  $\sigma_i^{k+1}$  is trained via Reinforcement Learning (or exact solving) to be a Best Response to the opponent’s meta-strategy:

$$\sigma_i^{k+1} \in \arg \max_{\sigma_i} \mathbb{E}_{\sigma_{-i} \sim \phi_{-i}} [u_i(\sigma_i, \sigma_{-i})] \quad (7)$$

In this work, we utilize an exact oracle that computes the optimal best response to the meta-strategy, isolating the performance of the meta-solver from the variance of reinforcement learning.

- (3) **Expansion:** The new policy is added to the population  $\Pi_i \leftarrow \Pi_i \cup \{\sigma_i^{k+1}\}$ , and the payoff tensor  $M$  is updated. In this work, we calculate the exact expected payoff value for each entry  $M_i^{j_1 \dots j_N}$ , thereby eliminating the stochastic noise associated with Monte Carlo sampling.

To quantify performance, we employ a distinct **evaluation-time meta-strategy solver** to compute a distribution over the current policy population; this distribution serves as the basis for the exploitability metric. It is critical [24] to distinguish this from the training process: the training-time MSS drives the *generation* of new policies (often prioritizing exploration), while the evaluation-time MSS ensures a robust *evaluation* of the population’s quality. We regard both the training-time and evaluation-time MSS as integral components of PSRO when employed as a single game-solver. Standard PSRO relies on fixed static meta-solvers for both training-time and eval-time (e.g., always using Nash or always using Uniform). However, the optimal meta-solver may change during training—for instance, encouraging exploration (Uniform) early on and robustness (Nash) later. We aim to discover a dynamic meta-solver schedule.

## 3 METHOD: AUTOMATING ALGORITHM DISCOVERY VIA ALPHAEVOLVE

We propose a framework to automate the design of multi-agent learning algorithms by shifting the design process from manual heuristics to automated discovery. We utilize **AlphaEvolve** [17], an evolutionary coding agent that leverages Large Language Models (LLMs) to evolve executable code. We apply this framework to two distinct paradigms: Regret Minimization (CFR) and Population based Training (PSRO).

### 3.1 The AlphaEvolve Framework

AlphaEvolve combines the creative code-generation capabilities of LLMs with the rigorous selection pressure of evolutionary algorithms. Unlike traditional genetic programming which relies on random syntactic mutations, AlphaEvolve uses an LLM to propose semantically meaningful modifications to a code file. The process operates as follows:

- **Population Initialization:** We initialize a population  $\mathcal{P}$  containing the standard implementations of the baseline algorithms (e.g., standard CFR code or Uniform PSRO code).
- **LLM-Driven Mutation:** In each generation, a parent algorithm  $A \in \mathcal{P}$  is selected based on fitness. Notice that AlphaEvolve supports *multi-objective* optimization: if multiple fitness metrics are defined, one of them (say  $f$ ) will be randomly selected, and  $A$  is then sampled in favor of high value of  $f$ . The source code of  $A$  is fed into an LLM (e.g., Gemini [7]) with a prompt instructing it to *"Modify the following code to improve fitness (reduce exploitability)."* The LLM generates a new candidate variant  $A'$  by applying several lines of code changes to  $A$ .
- **Automated Evaluation:** The candidate  $A'$  is executed on a set of proxy games (e.g., Kuhn Poker). An automated evaluator computes its fitness scores (final negative exploitability).
- **Evolutionary Selection:** Valid candidates are added to the population  $\mathcal{P}$ . This loop repeats, allowing the system to discover complex, non-intuitive optimizations.

### 3.2 Evolving Regret Minimization Code

To discover new variants of CFR, we expose the core update functions of the Counterfactual Regret Minimization algorithm to the AlphaEvolve agent. The search space is defined by the Python functions responsible for accumulating regret and updating the average policy. Specifically, we design the components with the primitives in Listing 1. This search space is expressive enough to encompass the entire family of known CFR variants as special cases. For instance, standard CFR uses eq (5) for `RegretAccumulator` and eq (6) for `PolicyFromRegretAccumulator`. These components are Python classes which allow for maintaining state across iterations, in contrast to pure functions.

**Listing 1: The Python CFR code skeleton used as the search space for AlphaEvolve. The highlighted methods `update_accumulate_regret`, `get_updated_current_policy`, and `update_accumulate_policy` represent the evolvable components of the CFR algorithm.**

```

1 class RegretAccumulator:
2     """A class that updates cumulative regret at an information set.
3         """
4     def update_accumulate_regret(self, info_state_node,
5                                 iteration_number, cfr_regrets):
6         """
7         Args:
8             info_state_node: Data structure with cumulative_regret and
9                             cumulative_policy.
10            iteration_number: Current CFR iteration.
11            cfr_regrets: Counterfactual regrets (not yet added to node).
12        Returns:
13            Updated cumulative regret dictionary for each action.
14        """
15        ...

```

```

15
16 class PolicyFromRegretAccumulator:
17     """A class that derives the current policy from regret."""
18
19     def get_updated_current_policy(self, info_state_node,
20         iteration_number, cfr_regrets, previous_policy):
21         """
22         Args:
23             info_state_node: Data structure with cumulative_regret.
24             iteration_number: Current CFR iteration.
25             cfr_regrets: Counterfactual regrets (already added to node).
26             previous_policy: Previous policy at this info set.
27         Returns:
28             Updated current policy dictionary.
29         """
30         ...
31
32 class PolicyAccumulator:
33     """A class that updates the average policy during tree traversal
34         ."""
35
36     def update_accumulate_policy(self, info_state_node,
37         iteration_number, info_state_policy, cfr_regrets,
38         reach_prob, counterfactual_reach_prob):
39         """
40         Args:
41             info_state_node: Data structure with cumulative_policy.
42             iteration_number: Current CFR iteration.
43             info_state_policy: Current policy at this info set.
44             cfr_regrets: Counterfactual regrets (already added to node).
45             reach_prob: Probability of reaching current history (player's
46                 contribution).
47             counterfactual_reach_prob: Probability of reaching current
48                 history (opponents' contribution).
49         Returns:
50             Updated cumulative policy dictionary.
51         """
52         ...

```

### 3.3 Evolving Meta-Strategy Solvers Code

For PSRO, we evolve TrainMetaStrategySolver for generating the mixed strategies used during the oracle training phase, and EvalMetaStrategySolver for computing a strategy profile to report metrics such as exploitability. The code skeleton we designed is shown in Listing 2. This interface supports the representation of all standard baselines as special cases: for example, standard double oracle algorithm solves for a Nash equilibrium in the get\_meta\_strategy method for both solver classes.

**Listing 2: The Python PSRO code skeleton used as the search space for AlphaEvolve. The highlighted methods TrainMetaStrategySolver and EvalMetaStrategySolver represent the evolvable components of the PSRO algorithm.**

```

1 class TrainMetaStrategySolver:
2     """Returns meta strategies to train against in PSRO."""
3
4     def get_meta_strategy(self, game, policy_sets, meta_games):
5         """Returns meta strategies to train against in policy-space
6             response oracles.
7
8         Args:
9             game: The pypspiel game object.
10            policy_sets: A list of lists of policies, one list per
11                player.
12            policy_sets[p][i] is player p's i-th policy. len(
13                policy_sets[p]) ==
14            meta_games[0].shape[p].
15            meta_games: A list of n-dimensional numpy arrays, one per
16                player. Each
17                array has shape (num_strats_p0, num_strats_p1, ...,
18                num_strats_pn-1) and

```

```

14            meta_games[p][i0, i1, ..., in-1] is the payoff of player p
15                when player k
16                chooses strategy ik.
17
18         Returns:
19             A list of mixed-strategies, one for each player. Each mixed
20                 strategy is
21                 a list of non-negative weights (not necessarily normalized).
22                 It is used
23                 to train best response against.
24         """
25         ...
26
27 class EvalMetaStrategySolver:
28     """Returns meta strategies for evaluation in PSRO."""
29
30     def get_meta_strategy(self, game, policy_sets, meta_games):
31         """Returns meta strategies for evaluation in policy-space
32             response oracles."""
33         ...

```

### 3.4 Optimization Objective

We manually select a set  $\mathcal{G}$  of training games for AlphaEvolve to compute  $|\mathcal{G}| + 1$  fitness scores. These are the negative exploitability  $-\text{Exploitability}(A(g)_K)$  of the final strategy profile after  $K$  iterations for each  $g \in \mathcal{G}$ , as well as their average

$$-\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \text{Exploitability}(A(g)_K).$$

Here  $A(g)_K$  denotes the strategy produced by algorithm  $A$  on game  $g$  at iteration  $K$ . The reported algorithms are selected based on their average scores.

## 4 EXPERIMENTAL EVALUATION

### 4.1 Experimental Setup

To test the robustness and generalizability of the algorithms we discover, we adopted a rigorous evaluation protocol involving two distinct sets of games. The algorithm's architecture and hyperparameters were developed and tuned on a *Training Set* of four games. For both CFR and PSRO discoveries, we choose this set to be 3-player Kuhn Poker, 2-player Leduc Poker, 4-card Goofspiel, and 5-sided Liars Dice. Subsequently, the fixed algorithm was evaluated on a *Test Set*. In the main body of this paper, we present results on four larger and more difficult games as a representative subset of test games: 4-player Kuhn Poker, 3-player Leduc Poker, 5-card Goofspiel, and 6-sided Liar's Dice. The results of a full-sweep of 11 games are provided in Appendix 7.3. We utilized the OpenSpiel [11] framework for all experiments. Our implementation of AlphaEvolve is backboned by Gemini 2.5 pro [7].

### 4.2 Experimental Evaluation: VAD-CFR<sup>1</sup>

We evaluate the efficacy of the evolved algorithm, VAD-CFR, against established game-theoretic baselines. We benchmarked VAD-CFR against a suite of state-of-the-art regret minimization algorithms: standard CFR [32], CFR+ [22], Linear CFR (LCFR), Discounted CFR (DCFR) [3], Predictive CFR+ (PCFR+) [8], Discounted Predictive CFR+ (DPCFR+) [28], and a Hyperparameter Schedule-powered

<sup>1</sup>Appendix 7.2 details an alternative variant (AOD-CFR) discovered on a different set of training games that also achieves robust results. Notably, this variant relies on more conventional mechanisms than VAD-CFR.

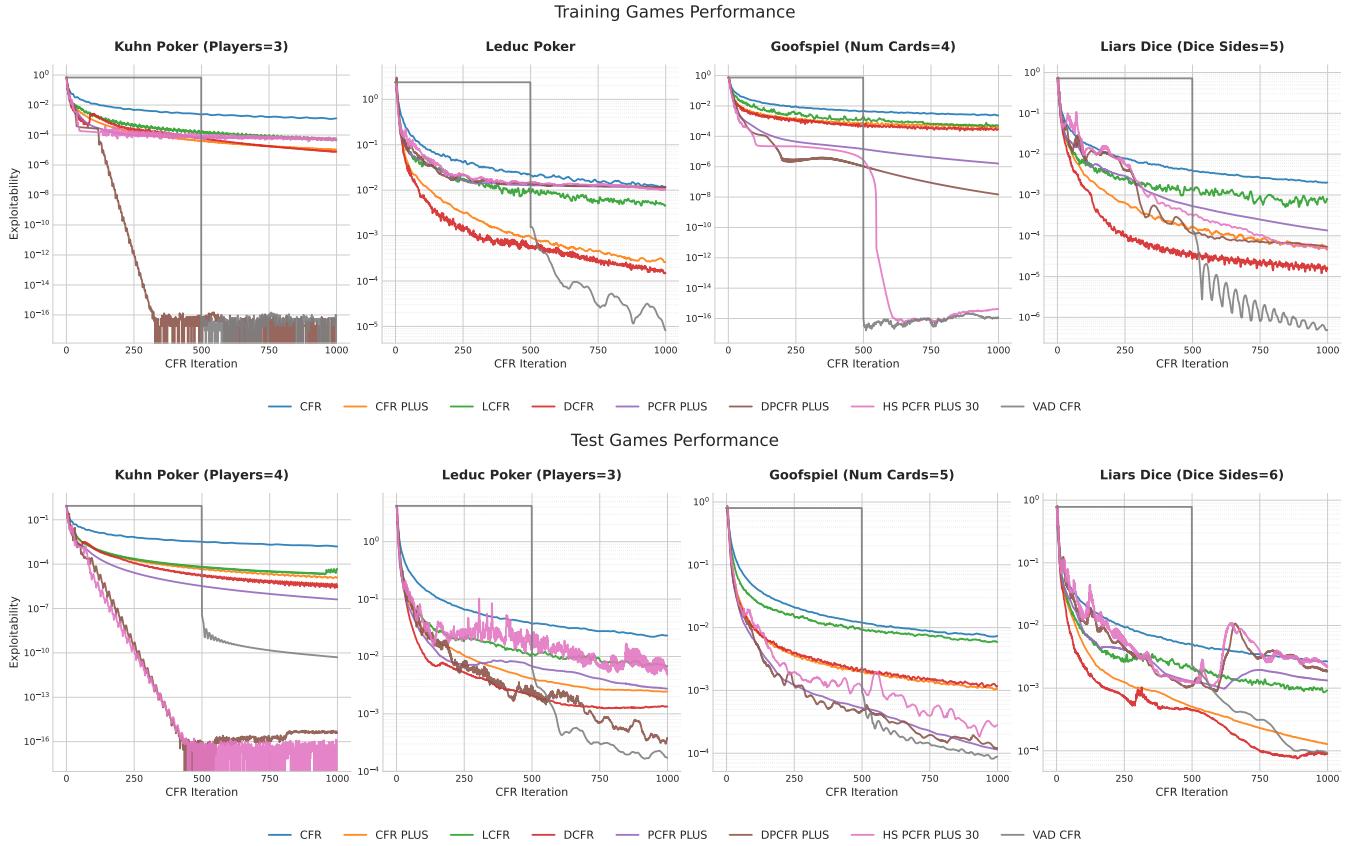


Figure 1: CFR variants performances.

PCFR variant HS-PCFR+(30) [31]. Performance was quantified using exploitability (measured on a logarithmic scale) over a fixed horizon of 1000 iterations, e.g.,  $K = 1000$ . For all domains the exploitability scores are computed exactly. We use CFR+ as the seed program. The prompt we use is shown in Listing 8 in the Appendix.

**4.2.1 Discovered Algorithmic Architecture.** The source code of the discovered VAD-CFR algorithm is provided in Listing 5 in Appendix 7.1. A high-level abstraction is shown in Listing 3. The evolutionary search discarded the standard linear averaging and static discounting of the CFR family in favor of three distinct, non-intuitive mechanisms:

- **Volatility-Adaptive Discounting (vs. Static Discounting):** Standard DCFR applies fixed discount factors ( $\alpha, \beta$ ) to cumulative regrets. In contrast, VAD-CFR makes these parameters reactive. It tracks the volatility of the learning process via an Exponential Weighted Moving Average (EWMA) of the instantaneous regret magnitude. When volatility is high (indicating the strategy is in flux), the algorithm dynamically increases discounting to "forget" the unstable history faster; when volatility drops, it retains more history for fine-tuning.
- **Asymmetric Instantaneous Boosting:** While DCFR applies asymmetry to accumulated history, VAD-CFR applies a novel asymmetry to the instantaneous update itself. Positive instantaneous

regrets (actions that are currently good) are boosted by a factor of 1.1. This asymmetry enables the immediate exploitation of beneficial deviations, eliminating the lag associated with accumulation.

- **Hard Warm-Start & Regret-Magnitude Weighting:** Perhaps the most interesting discovery is the policy accumulation schedule. Standard CFR variants begin averaging policies immediately at  $t = 1$ . In contrast, VAD-CFR enforces a hard warm-start<sup>2</sup>, postponing the start of policy averaging until iteration 500. The underlying regret accumulation process remains fully active during this phase, ensuring the agent continues to learn and update its strategy. Furthermore, once accumulation begins, policies are weighted not just linearly ( $t$ ), but by the magnitude of the instantaneous regret. This acts as a sophisticated filter, ensuring the final equilibrium approximation is constructed only from high-information iterations, preventing early-stage noise from polluting the solution quality. Warm starting mechanisms have been studied in a previous work [2].

**Listing 3: High-level abstraction of VAD-CFR.**

```

1 class RegretAccumulator:
2     """Volatility-Adaptive Discounting & Asymmetric Boosting."""
3     def update_accumulate_regret(self, info_state_node,
4                                 iteration_number, cfr_regrets):
5         # 1. Volatility & Adaptive Discount Calculation

```

<sup>2</sup>The LLM generated this threshold without knowing the fixed 1000-iteration evaluation horizon in its prompt context (Listing 8).

```

5     inst_mag = max((abs(r) for r in cfr_regrets.values()),
6                   default=0.0)
7     self.ewma = 0.1 * inst_mag + 0.9 * self.ewma
8     volatility = min(1.0, self.ewma / 2.0)
9     alpha = max(0.1, 1.5 - 0.5 * volatility)
10    beta = -0.1 - 0.5 * volatility
11    beta = min(alpha, beta)
12    t_plus_one = float(iteration_number) + 1.0
13    disc_pos = (t_plus_one ** alpha) / (t_plus_one ** alpha +
14    1.0)
15    disc_neg = (t_plus_one ** beta) / (t_plus_one ** beta +
16    1.0)
17    updated_regrets = {}
18    for a, r in cfr_regrets.items():
19        # 2. Asymmetric Instantaneous Boosting
20        r_boosted = r * 1.1 if r > 0 else r
21        # 3. Sign-Dependent History Discounting
22        prev_R = info_state_node.cumulative_regret.get(a, 0.0)
23        discount = disc_pos if prev_R >= 0 else disc_neg
24        # 4. Update with Negative Cap
25        new_R = (prev_R * discount) + r_boosted
26        updated_regrets[a] = max(-20.0, new_R)
27    return updated_regrets
28
29 class PolicyFromRegretAccumulator:
30     """Derives policy from a 'Future Projection' of regrets."""
31     def get_updated_current_policy(self, info_state_node,
32     iteration_number, cfr_regrets, previous_policy):
33         # 1. Consistency: Re-calculate exact adaptive params
34         inst_mag = max((abs(r) for r in cfr_regrets.values()),
35                       default=0.0)
36         volatility = min(1.0, inst_mag / 2.0)
37         # 2. Decaying Optimism Schedule
38         base_optimism = 1.0 / (1.0 + iteration_number / 100.0)
39         optimism = base_optimism * max(0.0, 1.0 - 0.5 * volatility)
40
41         policy = {}
42         sum_pos_regret = 0.0
43         for a in info_state_node.legal_actions:
44             # 3. Projected Regret Calculation
45             r_boosted = cfr_regrets.get(a, 0.0)
46             if r_boosted > 0: r_boosted *= 1.1
47             prev_R = info_state_node.cumulative_regret.get(a, 0.0)
48             discount = get_adaptive_discount_pos(prev_R) if prev_R
49             >= 0 else get_adaptive_discount_neg(prev_R)
50             proj_R = (prev_R * discount) + r_boosted + optimism
51             # 4. Non-Linear Probability Scaling
52             if proj_R > 0:
53                 scaled_r = proj_R ** 1.5
54                 policy[a] = scaled_r
55                 sum_pos_regret += scaled_r
56             else:
57                 policy[a] = 0.0
58         return normalize(policy) if sum_pos_regret > 0 else
59         uniform(policy)
60
61 class PolicyAccumulator:
62     """Warm-Start & Magnitude Weighting."""
63     def update_accumulate_policy(self, info_state_node,
64     iteration_number, info_state_policy, cfr_regrets,
65     reach_prob, counterfactual_reach_prob):
66         # 1. Hard Warm-Start
67         if iteration_number < 500:
68             return info_state_node.cumulative_policy
69         # 2. Adaptive Gamma (Polynomial Averaging)
70         inst_mag = max((abs(r) for r in cfr_regrets.values()),
71                       default=0.0)
72         volatility = min(1.0, inst_mag / 2.0)
73         gamma = min(4.0, 2.0 + 1.5 * volatility)
74         # 3. Multi-Factor Weighting
75         w_time = (float(iteration_number) + 1.0) ** gamma
76         w_mag = (1.0 + (inst_mag / 2.0)) ** 0.5
77         w_stable = 1.0 / (1.0 + inst_mag ** 1.5)
78         final_weight = w_time * w_mag * w_stable
79         # 4. Update Cumulative Policy
80         updated_policy = {}
81         for a, prob in info_state_policy.items():
82             prev_P = info_state_node.cumulative_policy.get(a, 0.0)

```

```

72         updated_policy[a] = prev_P + (final_weight *
73         reach_prob * prob)
74     return updated_policy

```

**4.2.2 Empirical Analysis: Training Domain.** As shown in Figure 1, VAD-CFR demonstrates superior convergence across the training set. For 3-player Kuhn Poker, VAD-CFR achieves significantly lower exploitability than all baselines. For Leduc Poker and 4-card Goofspiel, the algorithm maintains a steeper convergence slope compared to DPCFR+ and other state-of-the-art variants. In 5-sided Liars Dice, VAD-CFR exhibits robust performance, effectively managing the larger state space through its adaptive discounting and boosting mechanisms.

**4.2.3 Generalization Capabilities: Test Domain.** Results illustrated in Figure 1 highlights its generalization. In 3-player Leduc Poker, VAD-CFR reaches exploitability levels below  $10^{-3}$  while most baselines plateau at higher levels. In 6-sided Liar’s Dice, VAD-CFR continues to match established baselines like DCFR, suggesting that its evolved mechanisms for managing regret scaling and noise are highly effective across different game topologies.

Overall, VAD-CFR demonstrates efficient convergence rate and generalization across a broad variety of domains: as shown in Figure 3 in the Appendix, **VAD-CFR matches or surpasses previous state-of-the-art performances in 10 of the 11 games**, with 4-player Kuhn Poker as the only exception.

### 4.3 Experimental Evaluation: SHOR-PSRO

Next, we evaluated the performance of the evolved SHOR-PSRO algorithm, comparing its ability to reduce exploitability against standard meta-solver baselines. We use the exploitability at the  $K=100$ -th PSRO iteration as the metric. We employ an *exact best response oracle* via value iteration, where at each state it uniformly randomizes among actions with the same optimal values. We benchmarked SHOR-PSRO against standard established meta-solver baselines: Uniform, Nash equilibrium via linear program for 2-player games, AlphaRank [15], Projected Replicator Dynamics (PRD) [12], and Regret Matching (RM). We use Uniform as the initial program for both solver classes. The prompt we use is shown in Listing 9 in the Appendix.

**4.3.1 Discovered Algorithmic Architecture.** The source code of the discovered SHOR-PSRO algorithm is provided in Listing 6 in Appendix 7.1. A high-level abstraction is in Listing 4. The evolutionary search yielded a “Hybrid” meta-solver that constructs a meta-strategy  $\sigma$  by linearly blending two distinct components: Optimistic Regret Matching (ORM) and a Smoothed Best Pure Strategy.

• **Hybrid Blending Mechanism:** At every internal solver iteration, the meta-strategy is computed as:

$$\sigma_{\text{hybrid}} = (1 - \lambda) \cdot \sigma_{\text{ORM}} + \lambda \cdot \sigma_{\text{Softmax}} \quad (8)$$

where  $\sigma_{\text{ORM}}$  provides the stability of regret minimization, and  $\sigma_{\text{Softmax}}$  is a Boltzmann distribution over pure strategies that aggressively biases the solver toward high-reward modes. The blending factor  $\lambda$  controls the trade-off between these dynamics.

• **Dynamic Annealing Schedule:** Unlike standard PSRO which uses fixed solvers, SHOR-PSRO employs a dynamic schedule for the training-time solver. Over the course of the PSRO iterations:

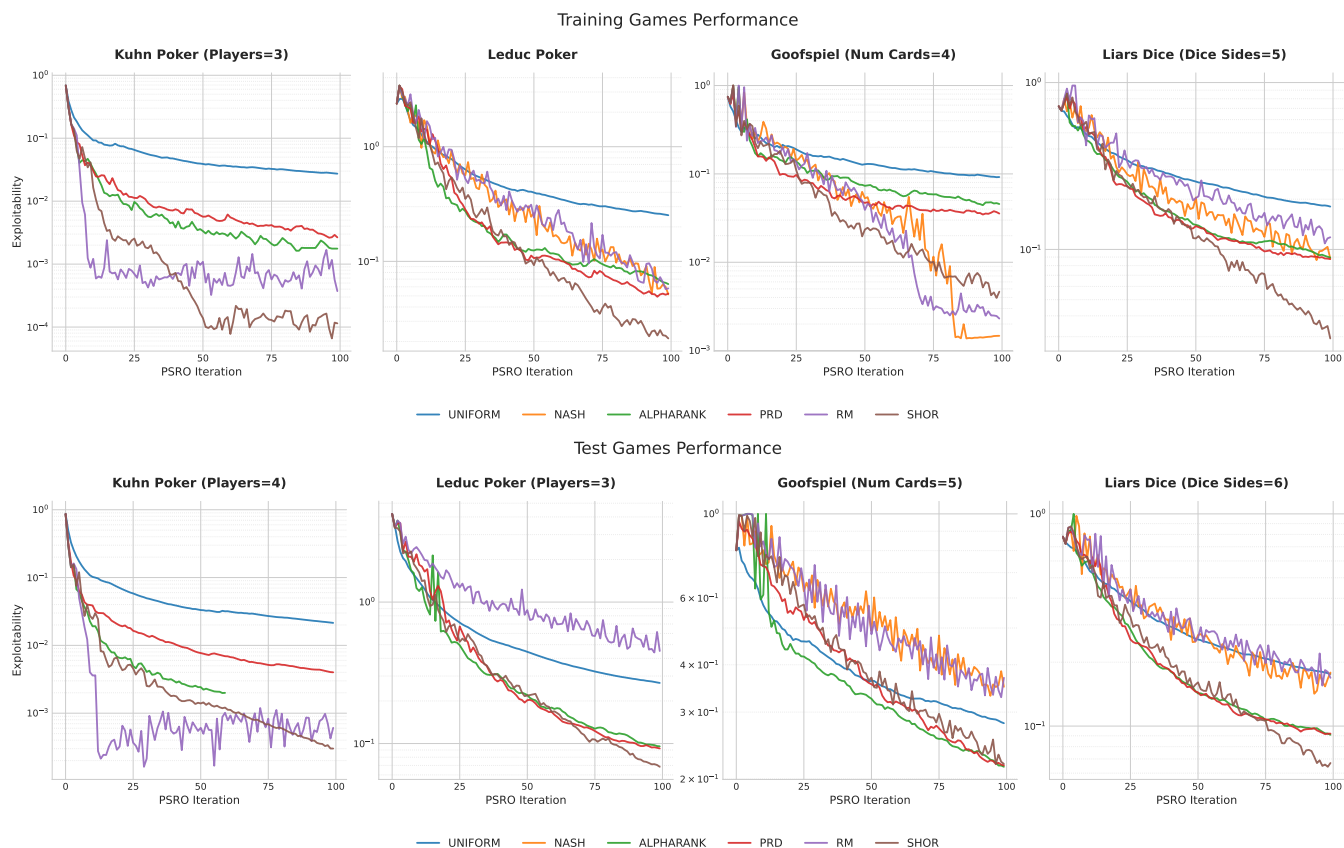


Figure 2: PSRO variants performances.

- The blending factor  $\lambda$  anneals from 0.3  $\rightarrow$  0.05, gradually shifting focus from greedy exploitation to robust equilibrium finding.
- An explicit **diversity bonus** is applied to the payoff gains, which also decays from 0.05  $\rightarrow$  0.001, ensuring early expansion of the game graph followed by late-stage refinement.

- **Training vs. Evaluation Asymmetry:** The search discovered distinct configurations for training and evaluation. The **Training Solver** utilizes the dynamic annealing schedule described above and returns the *average* strategy over internal iterations to ensure stability. In contrast, the **Evaluation Solver** employs a fixed, low blending factor ( $\lambda = 0.01$ ) and returns the *last-iterate* strategy. This decoupling allows the algorithm to explore safely during training while providing reactive, low-noise exploitability estimates during evaluation.

Listing 4: High-level abstraction of SHOR-PSRO.

```

1 class TrainMetaStrategySolver:
2     def get_meta_strategy(self, game, policy_sets, meta_games):
3         """Hybrid Training Solver with Full Parameter Annealing"""
4         # 1. Annealing Schedules (Exploration -> Exploitation)
5         prog = min(1.0, self.current_psro_iter / 75.0)
6         blend = 0.30 - (0.25 * prog) # Blend: 0.3 -> 0.05
7         div = 0.05 - (0.049 * prog) # Diversity: 0.05 -> 0.001
8         temp = 0.50 - (0.49 * prog) # Temp: 0.5 -> 0.01
9         # 2. Adaptive Solver Iterations
10        pop_size = len(meta_games[0])
11        n_iters = int(1000 + 20 * (pop_size - 1))

```

```

12        # 3. Hybrid Internal Loop (Training returns Average)
13        return _hybrid_solver(meta_games, n_iters, blend, div,
14                               temp, momentum=0.5, return_avg=True)
15
16 class EvalMetaStrategySolver:
17     def get_meta_strategy(self, game, policy_sets, meta_games):
18         """Evaluation Solver (Fixed Parameters, Last-Iterate)"""
19         # 1. Fixed Parameters for Pure Exploitation
20         blend, div, temp = 0.01, 0.0, 0.001
21         # 2. Aggressive Iteration Scale for Convergence
22         pop_size = len(meta_games[0])
23         n_iters = int(8000 + 50 * (pop_size - 1))
24         # 3. Hybrid Loop returns Last-Iterate for eval reactivity
25         return _hybrid_solver(meta_games, n_iters, blend, div,
26                               temp, momentum=0.2, return_avg=False)
27
28 def _hybrid_solver(meta_games, n_iters, blend, div, temp, momentum,
29                   return_avg):
30     sigma = uniform_strategy(meta_games)
31     sigma_avg = zeros_like(sigma)
32     for i in range(n_iters):
33         # A. Optimistic Regret Matching (ORM)
34         gains = compute_payoffs(meta_games, sigma) + div * (1 -
35                               sigma)
36         gains_norm = normalize_scale(apply_momentum(gains, beta=
37                               momentum))
38         sigma_orm = regret_matching(gains_norm)
39         # B. Smoothed Best Pure Strategy (Softmax)
40         exp_vals = compute_expected_values(meta_games, sigma)
41         sigma_pure = softmax(exp_vals, temperature=temp)
42         # C. Hybrid Blending
43         sigma = ((1 - blend) * sigma_orm) + (blend * sigma_pure)

```

```

40     if return_avg: sigma_avg += sigma
41     return normalize(sigma_avg) if return_avg else sigma

```

**4.3.2 Empirical Analysis: Training Domain.** As shown in Figure 2, SHOR-PSRO demonstrates superior convergence speed and stability across the training set. In simpler domains like Kuhn Poker, SHOR-PSRO achieves exploitability levels ( $< 10^{-3}$ ) significantly faster than PRD or RM. This acceleration is likely attributable to the "Hybrid Blending" mechanism, which allows the solver to leverage the stability of regret minimization while aggressively exploiting high-reward modes via the smoothed softmax component.

**4.3.3 Generalization Capabilities: Test Domain.** To assess the robustness of the evolved mechanism, we evaluated SHOR-PSRO on a Test Set comprising larger and more complex game variants. The results, illustrated in Figure 2, highlight the algorithm’s generalization capabilities. In 3-player Leduc Poker, the meta-game landscape becomes significantly more chaotic due to multi-agent dynamics. Despite this, SHOR-PSRO consistently matches or outperforms the best-performing baselines. In the most demanding test case, Liar’s Dice (6 sides), SHOR-PSRO demonstrates a clear advantage. While static solvers struggle with the expanded branching factor, the hybrid solver’s ability to bias towards "smoothed best pure strategies" allows it to navigate the expanded game graph efficiently. The results suggest that the evolved evaluation-time asymmetry—using a fixed, low-noise configuration for measurement—ensures that the exploitability metric accurately reflects the population’s growing strength without being obscured by exploration noise.

Overall, SHOR-PSRO proves to be a highly robust meta-solver, capable of generalizing to unseen game topologies without requiring manual re-tuning of its annealing schedules. As shown in Figure 4 in the Appendix, **SHOR-PSRO matches or surpasses state-of-the-art performances in 8 of the 11 games.**

## 5 RELATED WORK

Our work sits at the intersection of game-theoretic multi-agent learning, meta-learning, and automated algorithmic discovery. Here, we review the manual heuristic design that characterizes the current state-of-the-art and contrast it with our automated approach.

Significant research effort has been dedicated to improving the convergence speed of CFR through specific weighting schemes and regret targets. Notable variants include CFR+ [22], which introduced floor bounding and linear averaging; Discounted CFR (DCFR) [3], which applies fixed discount factors to historical regrets; and Predictive CFR+ (PCFR+) [8], which utilizes predictive Blackwell approachability. However, these improvements were largely derived through human intuition and trial-and-error to navigate the vast combinatorial space of potential update rules. In contrast to strategy-based warm starting [2], which relies on discrete, heuristic-driven resets, VAD-CFR employs a static warmup schedule with continuous volatility-based discounting to regulate regret retention.

For larger games where traversing the full tree is infeasible, Policy Space Response Oracles (PSRO) [1] generalizes the Double Oracle algorithm [13] by iteratively expanding a population of policies via exact or reinforcement learning best responses. While PSRO rests on solid theoretical ground [30], its practical application faces challenges in terms of convergence speed and population quality. Recent surveys highlight the diversity of PSRO variants [1],

but like CFR, the design of meta-solver schedules has traditionally been static or manually tuned rather than dynamically evolved.

The research of automating machine learning algorithm design has been evolving for both neural-based approaches and symbolic-based approaches. Meta-learning approaches such as meta-RL [18, 29] and meta-learning optimizers [14] have parameterized update rules using neural networks to optimize learning dynamics. For symbolic ML discovery, a foundational work in this domain is AutoML-Zero [19], which demonstrated that complete machine learning algorithms could be evolved from scratch using basic mathematical operations. A following work also applies program search for discovering optimizers [5]. The authors of [6] also studied discovering novel symbolic reinforcement learning algorithms.

Within the specific domain of multi-agent learning, there has been prior work on automating algorithm design [9, 26, 27]. More directly related, one can meta-learn the regret minimization algorithm directly via a CFR-like no-regret framework from examples [21]. However, these approaches often faced limitations: they either operated within a relatively constrained search space or relied on neural parameterizations that hindered interpretability. Our work builds directly on AlphaEvolve [17], utilizing Large Language Models (LLMs) to perform semantic mutation on interpretable code, bridging the gap between expressive neural meta-learning and symbolic discovery. This approach has already been shown great success in the field of math [10] and combinatorial algorithms [16].

## 6 CONCLUSION

In this work, we proposed a methodology for using LLM-driven evolution to automate the design of multi-agent learning algorithms, effectively shifting the paradigm from manual heuristic tuning to symbolic code evolution. By leveraging AlphaEvolve, we demonstrated that algorithm design can be treated as a search problem within a combinatorial space of symbolic operations, allowing for the discovery of complex, non-intuitive optimization mechanisms.

First, we identified **VAD-CFR**, which breaks from the tradition of static discount schedules used in DCFR. By introducing volatility-adaptive parameters, instantaneous regret boosting, and a warm-start for policy averaging, VAD-CFR demonstrates that filtering out early-stage noise is as critical as optimizing late-stage convergence. These mechanisms significantly outperform state-of-the-art baselines in domains where exact exploitability is computable.

Second, we addressed the efficiency challenges of Policy Space Response Oracles (PSRO) by evolving **SHOR-PSRO**. This algorithm utilizes a hybrid meta-solver that blends Optimistic Regret Matching with a smoothed, temperature-controlled best response. By automating the annealing of this blending factor, SHOR-PSRO effectively manages the transition from population diversity to equilibrium refinement, yielding solvers that are significantly more effective for equilibrium finding in large-scale scenarios.

Our results suggest that the automated discovery of algorithmic asymmetries—specifically those that manage regret scaling and dynamic mixing schedules—can yield solvers that are elusive to human intuition but highly effective in practice. Future work will explore the application of this evolutionary framework to fully deep reinforcement learning agents and the discovery of cooperative mechanisms in general-sum games.

## REFERENCES

- [1] Ariyan Bighashdel, Yongzhao Wang, Stephen McAleer, Rahul Savani, and Frans A Oliehoek. 2024. Policy Space Response Oracles: A Survey. In *Thirty-Third International Joint Conference on Artificial Intelligence (IJCAI-24) Survey Track*.
- [2] Noam Brown and Tuomas Sandholm. 2016. Strategy-based warm starting for regret minimization in games. In *Thirtieth AAAI Conference on Artificial Intelligence*, Vol. 30.
- [3] Noam Brown and Tuomas Sandholm. 2019. Solving imperfect-information games via discounted regret minimization. In *Thirty-Third AAAI Conference on Artificial Intelligence*, Vol. 33. 1829–1836.
- [4] Noam Brown and Tuomas Sandholm. 2019. Superhuman AI for multiplayer poker. *Science* 365, 6456 (2019), 885–890. <https://doi.org/10.1126/science.aay2400>
- [5] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, et al. 2023. Symbolic discovery of optimization algorithms. *Thirty-Seventh International Conference on neural information processing systems*, 49205–49233.
- [6] John D Co-Reyes, Yingjie Miao, Daiyi Peng, Esteban Real, Quoc V Le, Sergey Levine, Honglak Lee, and Aleksandra Faust. [n.d.]. Evolving Reinforcement Learning Algorithms. In *Ninth International Conference on Learning Representations*.
- [7] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261* (2025).
- [8] Gabriele Farina, Christian Kroer, and Tuomas Sandholm. 2021. Faster game solving via predictive blackwell approachability: Connecting regret matching and mirror descent. In *Thirty-Fifth AAAI Conference on Artificial Intelligence*, Vol. 35. 5363–5371.
- [9] Xidong Feng, Oliver Slumbers, Ziyu Wan, Bo Liu, Stephen Marcus McAleer, Ying Wen, Jun Wang, and Yaodong Yang. [n.d.]. Neural Auto-Curricula in Two-Player Zero-Sum Games. In *Thirty-Fifth International Conference on Neural Information Processing Systems*.
- [10] Bogdan Georgiev, Javier Gómez-Serrano, Terence Tao, and Adam Zsolt Wagner. 2025. Mathematical exploration and discovery at scale. *arXiv preprint arXiv:2511.02864* (2025).
- [11] Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay, Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, et al. 2019. OpenSpiel: A framework for reinforcement learning in games. *arXiv preprint arXiv:1908.09453* (2019).
- [12] Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, and Thore Graepel. 2017. A Unified Game-Theoretic Approach to Multiagent Reinforcement Learning. In *Thirty-First International Conference on Neural Information Processing Systems*. 4190–4203.
- [13] H Brendan McMahan, Geoffrey J Gordon, and Avrim Blum. 2003. Planning in the presence of cost functions controlled by an adversary. In *Twentieth International Conference on Machine Learning*. 536–543.
- [14] Luke Metz, Niru Maheswaranathan, Jeremy Nixon, C Daniel Freeman, and Jascha Sohl-Dickstein. 2019. Understanding and Correcting Pathologies in the Training of Learned Optimizers. In *Thirty-Sixth International Conference on Machine Learning*. 4556–4565.
- [15] Paul Muller, Shayegan Omidshafiei, Mark Rowland, Karl Tuyls, Julien Perolat, Siqi Liu, Daniel Hennes, Luke Marris, Marc Lanctot, Edward Hughes, et al. [n.d.]. A Generalized Training Approach for Multiagent Learning. In *Ninth International Conference on Learning Representations*.
- [16] Ansh Nagda, Prabhakar Raghavan, and Abhradeep Thakurta. 2025. Reinforced Generation of Combinatorial Structures: Hardness of Approximation. *arXiv preprint arXiv:2509.18057* (2025).
- [17] Alexander Novikov, Ngàn Vù, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, et al. 2025. AlphaEvolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint arXiv:2506.13131* (2025).
- [18] Junhyuk Oh, Greg Farquhar, Iurii Kemaev, Dan A Calian, Matteo Hessel, Luisa Zintgraf, Satinder Singh, Hado Van Hasselt, and David Silver. 2025. Discovering state-of-the-art reinforcement learning algorithms. *Nature* (2025), 1–2.
- [19] Esteban Real, Chen Liang, David So, and Quoc V Le. 2020. AutoML-Zero: Evolving Machine Learning Algorithms From Scratch. In *Thirty-Seventh International Conference on Machine Learning*. 8007–8019.
- [20] Yoav Shoham and Kevin Leyton-Brown. 2008. *Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations*. Cambridge University Press, New York, NY, USA.
- [21] David Sychrovsk’y, Michal Šustr, Elnaz Davoodi, Michael Bowling, Marc Lanctot, and Martin Schmid. [n.d.]. Learning to Not Regret. In *Thirty-Eighth AAAI Conference on Artificial Intelligence*.
- [22] Oskari Tammelin. 2014. Solving large imperfect information games using CFR+. *arXiv preprint arXiv:1407.5042* (2014).
- [23] Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* 575, 7782 (2019), 350–354.
- [24] Yongzhao Wang, Qiurui Ma, and Michael P Wellman. 2022. Evaluating Strategy Exploration in Empirical Game-Theoretic Analysis. In *Twenty-First International Conference on Autonomous Agents and Multiagent Systems*. 1346–1354.
- [25] Michael P Wellman, Karl Tuyls, and Amy Greenwald. 2025. Empirical game theoretic analysis: A survey. *Journal of Artificial Intelligence Research* 82 (2025), 1017–1076.
- [26] Hang Xu, Kai Li, Haobo Fu, Qiang Fu, and Junliang Xing. 2022. Autocfr: Learning to design counterfactual regret minimization algorithms. In *Thirty-Sixth AAAI Conference on Artificial Intelligence*. 5244–5251.
- [27] Hang Xu, Kai Li, Haobo Fu, Qiang Fu, Junliang Xing, and Jian Cheng. 2024. Dynamic discounted counterfactual regret minimization. In *Twelfth International Conference on Learning Representations*.
- [28] Hang Xu, Kai Li, Bingyun Liu, Haobo Fu, Qiang Fu, Junliang Xing, and Jian Cheng. 2024. Minimizing weighted counterfactual regret with optimistic online mirror descent. In *Thirty-Third International Joint Conference on Artificial Intelligence*. 5272–5280.
- [29] Zhongwen Xu, Hado van Hasselt, and David Silver. 2018. Meta-Gradient Reinforcement Learning. In *Thirty-Second International Conference on Neural Information Processing Systems*, Vol. 31.
- [30] Brian Hu Zhang and Tuomas Sandholm. 2024. Exponential lower bounds on the double oracle algorithm in zero-sum games. In *Thirty-Third International Joint Conference on Artificial Intelligence*. 3032–3039.
- [31] Naifeng Zhang, Stephen McAleer, and Tuomas Sandholm. 2026. Faster game solving via hyperparameter schedules. In *Fortieth AAAI Conference on Artificial Intelligence*.
- [32] Martin Zinkevich, Michael Johanson, Michael Bowling, and Carmelo Piccione. 2007. Regret Minimization in Games with Incomplete Information. In *Twenty-First International Conference on Neural Information Processing Systems*. 1729–1736.

## 7 APPENDIX

### 7.1 Source code of discovered algorithms

Listing 5: VAD-CFR

```
1 class RegretAccumulator:
2     """A class that updates cumulative regret using Adaptive Discounting
3     with separate discounting for positive and negative regrets, and instantaneous regret boosting.
4     """
5
6     @staticmethod
7     def _calculate_adaptive_params(
8         iteration_number,
9         cfr_regrets,
10        base_alpha,
11        base_beta,
12        volatility_sensitivity,
13        max_expected_instantaneous_regret,
14        ewma_decay_factor,
15        current_ewma_magnitude,
16    ):
17        """Calculates adaptive discounting parameters for a given iteration.
18        This static method centralizes the logic for EWMA volatility, adaptive
19        alpha/beta, and discount factors to ensure consistency across components.
20        """
21        t_plus_one = float(iteration_number + 1)
22
23        instantaneous_regret_magnitude = max(
24            (abs(r) for r in cfr_regrets.values()), default=0.0
25        )
26
27        if iteration_number == 0:
28            projected_ewma = instantaneous_regret_magnitude
29        else:
30            projected_ewma = (
31                ewma_decay_factor * instantaneous_regret_magnitude +
32                (1.0 - ewma_decay_factor) * current_ewma_magnitude
33            )
34
35        if max_expected_instantaneous_regret > 0:
36            normalized_volatility = min(1.0, projected_ewma / max_expected_instantaneous_regret)
37        else:
38            normalized_volatility = 0.0
39
40        effective_alpha = max(0.1, base_alpha - volatility_sensitivity * normalized_volatility)
41        effective_beta = base_beta - volatility_sensitivity * normalized_volatility
42        effective_beta = min(effective_alpha, effective_beta)
43
44        discount_factor_positive = (t_plus_one**effective_alpha) / (t_plus_one**effective_alpha + 1.0)
45        discount_factor_negative = (t_plus_one**effective_beta) / (t_plus_one**effective_beta + 1.0)
46
47        return projected_ewma, normalized_volatility, discount_factor_positive, discount_factor_negative
48
49    def __init__(self, base_alpha=1.5, base_beta=-0.1, volatility_sensitivity=0.5,
50                max_expected_instantaneous_regret=2.0, instantaneous_regret_boost_factor=1.1,
51                ewma_decay_factor=0.1, negative_regret_cap=-20.0):
52        """Initializes the regret accumulator with adaptive discounting parameters.
53
54        Args:
55            base_alpha: The baseline exponent for discounting positive cumulative regrets.
56            base_beta: The baseline exponent for discounting negative cumulative regrets.
57            volatility_sensitivity: Controls how strongly the instantaneous regret
58                magnitude influences the adaptive alpha/beta. A higher value means the
59                exponents will be more reduced by high volatility.
60            max_expected_instantaneous_regret: An estimate of the maximum possible
61                instantaneous regret magnitude, used for normalizing the volatility.
62            instantaneous_regret_boost_factor: Boost factor for positive instantaneous regrets.
63                A factor > 1.0 makes the algorithm more reactive to current good actions.
64            ewma_decay_factor: Decay factor for the EWMA of instantaneous regret magnitude.
65            negative_regret_cap: The minimum value for cumulative regret, to prevent
66                regret lock-in and improve adaptability.
67        """
68        self._base_alpha = base_alpha
69        self._base_beta = base_beta
70        self._volatility_sensitivity = volatility_sensitivity
71        self._max_expected_instantaneous_regret = max_expected_instantaneous_regret
72        self._instantaneous_regret_boost_factor = instantaneous_regret_boost_factor
```

```

73 self._ewma_decay_factor = ewma_decay_factor
74 self._negative_regret_cap = negative_regret_cap
75 self._ewma_instantaneous_regret_magnitude = 0.0
76
77 def update_accumulate_regret(
78     self, info_state_node, iteration_number, cfr_regrets
79 ):
80     """Updates cumulative regret for each action at an information set using ADCFRP.
81     Cumulative regrets are now signed.
82
83     Args:
84         info_state_node: a data structure corresponding to an information set with
85             cumulative_regret and cumulative_policy stored.
86         iteration_number: the current CFR iteration (0-indexed).
87         cfr_regrets: the instantaneous counterfactual regrets of the current policy at the
88             current information set. cfr_regrets haven't been added to info_state_node.cumulative_regret.
89     Returns:
90         updated cumulative regret for each action at the current information set (signed).
91     """
92     # Centralize adaptive parameter calculation to ensure consistency and reduce redundancy.
93     (
94         self._ewma_instantaneous_regret_magnitude,
95         _, # normalized_volatility is not used here
96         discount_factor_positive,
97         discount_factor_negative,
98     ) = RegretAccumulator._calculate_adaptive_params(
99         iteration_number=iteration_number,
100         cfr_regrets=cfr_regrets,
101         base_alpha=self._base_alpha,
102         base_beta=self._base_beta,
103         volatility_sensitivity=self._volatility_sensitivity,
104         max_expected_instantaneous_regret=self._max_expected_instantaneous_regret,
105         ewma_decay_factor=self._ewma_decay_factor,
106         current_ewma_magnitude=self._ewma_instantaneous_regret_magnitude,
107     )
108
109     updated_cumulative_regret = {}
110     for action in cfr_regrets:
111         old_regret = info_state_node.cumulative_regret[action]
112
113         instantaneous_regret_component = cfr_regrets[action]
114         if instantaneous_regret_component > 0:
115             instantaneous_regret_component *= self._instantaneous_regret_boost_factor
116
117         # Apply different discount factors based on the sign of the old regret.
118         if old_regret >= 0:
119             discounted_old_regret = discount_factor_positive * old_regret
120         else:
121             discounted_old_regret = discount_factor_negative * old_regret
122
123         new_regret = discounted_old_regret + instantaneous_regret_component
124
125         # Cap the negative regret to prevent lock-in and improve adaptability.
126         new_regret = max(self._negative_regret_cap, new_regret)
127
128         # Crucially, regrets are NOT clipped to zero here. They can be negative.
129         updated_cumulative_regret[action] = new_regret
130
131     return updated_cumulative_regret
132
133
134 class PolicyFromRegretAccumulator:
135     """A class that obtains a current policy from a consistent optimistic projection of regrets.
136     It aligns the policy generation with the adaptive, asymmetric, and boosted logic from RegretAccumulator.
137     """
138
139     def __init__(self,
140                 initial_optimism_factor=1.0,
141                 optimism_decay_factor=100.0,
142                 positive_policy_exponent=1.5,
143                 base_alpha=1.5,
144                 base_beta=-0.1,
145                 volatility_sensitivity=0.5,
146                 max_expected_instantaneous_regret=2.0,
147                 instantaneous_regret_boost_factor=1.1,
148                 ewma_decay_factor=0.1):
149         """Initializes the PolicyFromRegretAccumulator with parameters consistent with RegretAccumulator.
150

```

```

151 Args:
152     initial_optimism_factor: The initial weighting for the instantaneous regret component in the projection.
153     optimism_decay_factor: Controls how quickly the optimism weight decays.
154     positive_policy_exponent: Exponent for non-linear scaling of positive regrets.
155     base_alpha: The baseline exponent for discounting positive regrets.
156     base_beta: The baseline exponent for discounting negative regrets.
157     volatility_sensitivity: Controls influence of volatility on both discounting
158         and optimism dampening.
159     max_expected_instantaneous_regret: Used for normalizing volatility.
160     instantaneous_regret_boost_factor: Boost factor for positive instantaneous regrets.
161     ewma_decay_factor: Decay factor for the EWMA of volatility.
162     """
163     self._initial_optimism_factor = initial_optimism_factor
164     self._optimism_decay_factor = optimism_decay_factor
165     self._positive_policy_exponent = positive_policy_exponent
166     self._base_alpha = base_alpha
167     self._base_beta = base_beta
168     self._volatility_sensitivity = volatility_sensitivity
169     self._max_expected_instantaneous_regret = max_expected_instantaneous_regret
170     self._instantaneous_regret_boost_factor = instantaneous_regret_boost_factor
171     self._ewma_decay_factor = ewma_decay_factor
172     self._ewma_instantaneous_regret_magnitude = 0.0
173
174 def get_updated_current_policy(self, info_state_node, iteration_number, cfr_regrets, previous_policy):
175     """Obtains the current policy using a projection that is consistent with the RegretAccumulator's update rule.
176
177     This method creates a tighter feedback loop by basing the current policy
178     on a projection of what the regrets will be *after* the current
179     iteration's update. This projection uses the same adaptive, asymmetric
180     discounting and boosting logic as the main regret accumulation step.
181     """
182     # Centralize adaptive parameter calculation, ensuring consistency with RegretAccumulator.
183     (
184         self._ewma_instantaneous_regret_magnitude,
185         normalized_volatility,
186         discount_factor_positive,
187         discount_factor_negative,
188     ) = RegretAccumulator._calculate_adaptive_params(
189         iteration_number=iteration_number,
190         cfr_regrets=cfr_regrets,
191         base_alpha=self._base_alpha,
192         base_beta=self._base_beta,
193         volatility_sensitivity=self._volatility_sensitivity,
194         max_expected_instantaneous_regret=self._max_expected_instantaneous_regret,
195         ewma_decay_factor=self._ewma_decay_factor,
196         current_ewma_magnitude=self._ewma_instantaneous_regret_magnitude,
197     )
198
199     base_optimism = self._initial_optimism_factor / (1.0 + float(iteration_number) / self._optimism_decay_factor)
200     # Dampen optimism during volatile periods to increase stability.
201     optimism_dampening_factor = max(0.0, 1.0 - self._volatility_sensitivity * normalized_volatility)
202     optimism_strength = base_optimism * optimism_dampening_factor
203
204     action_to_projected_regret = {}
205     for action in info_state_node.legal_actions:
206         old_cumulative_regret = info_state_node.cumulative_regret.get(action, 0.0)
207         instantaneous_regret = cfr_regrets.get(action, 0.0)
208
209         instantaneous_regret_component = instantaneous_regret
210         if instantaneous_regret_component > 0:
211             instantaneous_regret_component *= self._instantaneous_regret_boost_factor
212
213         if old_cumulative_regret >= 0:
214             discounted_old_regret = discount_factor_positive * old_cumulative_regret
215         else:
216             discounted_old_regret = discount_factor_negative * old_cumulative_regret
217
218         projected_regret = discounted_old_regret + optimism_strength * instantaneous_regret_component
219         action_to_projected_regret[action] = projected_regret
220
221     positive_scaled_projected_regrets = {
222         action: (max(0.0, regret) ** self._positive_policy_exponent)
223         for action, regret in action_to_projected_regret.items()
224     }
225
226     sum_positive_scaled_projected_regrets = sum(positive_scaled_projected_regrets.values())
227
228     info_state_policy = {}

```

```

229 if sum_positive_scaled_projected_regrets > 0:
230     for action, scaled_regret in positive_scaled_projected_regrets.items():
231         info_state_policy[action] = scaled_regret / sum_positive_scaled_projected_regrets
232 else:
233     num_legal_actions = len(info_state_node.legal_actions)
234     for action in info_state_node.legal_actions:
235         info_state_policy[action] = 1.0 / num_legal_actions
236
237 return info_state_policy
238
239
240 class PolicyAccumulator:
241     """A class that updates cumulative policy using regret-informed weighted
242     averaging with a warmup period."""
243
244     def __init__(self, base_gamma=2.0, gamma_max=4.0,
245                 gamma_volatility_sensitivity=1.5, warmup_iterations=500,
246                 stability_exponent=1.5, max_expected_instantaneous_regret=2.0,
247                 regret_magnitude_weighting_exponent=0.5): # New parameter
248         """Initializes the PolicyAccumulator with adaptive gamma parameters and regret-magnitude weighting.
249
250     Args:
251         base_gamma: The baseline exponent for polynomial weighting of policies.
252         gamma_max: The maximum value the adaptive gamma can reach.
253         gamma_volatility_sensitivity: Controls how strongly volatility influences gamma.
254         warmup_iterations: Number of initial iterations to skip for policy averaging.
255         stability_exponent: Exponent for the stability factor based on regret magnitude.
256         max_expected_instantaneous_regret: Normalization factor for instantaneous regret magnitude.
257         regret_magnitude_weighting_exponent: Exponent for up-weighting policies based on
258         the absolute magnitude of instantaneous regrets. Higher values give more
259         emphasis to policies from iterations with large regrets.
260
261         """
262         self._base_gamma = base_gamma
263         self._gamma_max = gamma_max
264         self._gamma_volatility_sensitivity = gamma_volatility_sensitivity
265         self._warmup_iterations = warmup_iterations
266         self._stability_exponent = stability_exponent
267         self._max_expected_instantaneous_regret = max_expected_instantaneous_regret
268         self._regret_magnitude_weighting_exponent = regret_magnitude_weighting_exponent # Stored
269
270     def update_accumulate_policy(
271         self,
272         info_state_node,
273         iteration_number,
274         info_state_policy,
275         cfr_regrets,
276         reach_prob,
277         counterfactual_reach_prob,
278     ):
279         """Updates cumulative policy using delayed, regret-informed, and regret-magnitude weighted averaging.
280
281         """
282         if iteration_number < self._warmup_iterations:
283             return info_state_node.cumulative_policy
284
285         # Calculate instantaneous regret magnitude (L-infinity norm) for this iteration
286         instantaneous_regret_magnitude = max(
287             (abs(r) for r in cfr_regrets.values()), default=0.0
288         )
289
290         # Normalize volatility using the shared parameter
291         if self._max_expected_instantaneous_regret > 0:
292             normalized_volatility = min(1.0, instantaneous_regret_magnitude / self._max_expected_instantaneous_regret)
293         else:
294             normalized_volatility = 0.0
295
296         # Adapt gamma based on volatility: higher volatility -> higher gamma
297         effective_gamma = self._base_gamma + self._gamma_volatility_sensitivity * normalized_volatility
298         effective_gamma = min(self._gamma_max, effective_gamma)
299
300         # Standard polynomial weighting gives more weight to later iterations, now with adaptive gamma.
301         temporal_weight = (float(iteration_number) + 1.0) ** effective_gamma
302
303         # Calculate a stability factor from the L-infinity norm of instantaneous regrets.
304         # Higher regret magnitude -> lower stability factor, using L-infinity norm for consistency.
305         regret_stability_factor = 1.0 / (1.0 + instantaneous_regret_magnitude**self._stability_exponent)
306
307         # NEW: Regret Magnitude Weighting Factor
308         # Policies from iterations with higher regret magnitude contribute more to the average.

```

```

307 # This factor boosts the weight, with higher values of the exponent giving more emphasis.
308 # Ensure it's never zero to avoid division by zero or completely nullifying weight.
309 regret_magnitude_factor = (
310     1.0 + (instantaneous_regret_magnitude / self._max_expected_instantaneous_regret)
311 ) ** self._regret_magnitude_weighting_exponent
312 regret_magnitude_factor = max(0.1, regret_magnitude_factor) # Ensure minimum value to avoid zero weight if normalization results in very
    small number
313
314 # The final weight combines the temporal, stability, and regret-magnitude-based components.
315 weight = temporal_weight * regret_stability_factor * regret_magnitude_factor
316
317 return {
318     action: (
319         info_state_node.cumulative_policy[action]
320         + weight * reach_prob * info_state_policy[action]
321     )
322     for action in info_state_policy
323 }

```

### Listing 6: SHOR-PSRO

```

1 import numpy as np
2
3
4 def _smoothed_best_pure_strategy(payoff_vec, temperature=1.0):
5     """Computes a smoothed distribution biased towards the best pure strategy.
6     The softmax function ensures that strategies with higher payoffs are
7     given higher probability, with 'temperature' controlling the sharpness
8     of the distribution. A lower temperature makes the distribution more
9     concentrated on the best strategy, while a higher temperature
10    leads to a more uniform distribution.
11    """
12    # Subtract max payoff for numerical stability (standard softmax trick)
13    stable_payoffs = payoff_vec - np.max(payoff_vec)
14    exp_payoffs = np.exp(stable_payoffs / temperature)
15
16    sum_exp_payoffs = np.sum(exp_payoffs)
17    if sum_exp_payoffs > 1e-12: # Avoid division by zero
18        return exp_payoffs / sum_exp_payoffs
19    else:
20        # Fallback to uniform distribution if all exponentiated payoffs are
21        # effectively zero (e.g., due to very low temperature and negative payoffs,
22        # or all payoffs being identical after stabilization).
23        return np.ones_like(payoff_vec) / len(payoff_vec)
24
25 def _hybrid_orm_solver(meta_games, iterations,
26                       blending_factor=0.0,
27                       temperature=0.1,
28                       momentum_beta=0.0,
29                       gain_normalization=True,
30                       diversity_bonus_coeff=0.0,
31                       return_average_strategy=True): # New: Flag to return average or last-iterate strategy
32     """Computes meta-strategies using Optimistic Regret Matching+ enhanced with
33     optimistic updates, gain normalization, and a diversity bonus, then blended
34     with a smoothed best pure strategy.
35
36     This solver combines the stability and convergence properties of Optimistic
37     Regret Matching+ (ORM+) with an explicit pull towards highly rewarding
38     pure strategies, smoothed by a temperature-controlled softmax. This hybrid
39     approach aims to leverage ORM+'s ability to find mixed equilibria while
40     also quickly identifying and exploring strong pure-strategy modes in the
41     meta-game, thereby potentially accelerating the discovery of low-exploitable
42     policies in PSRO. The blending factor controls the trade-off between
43     these two dynamics.
44
45     Args:
46     meta_games: A list of n-dimensional numpy arrays, one per player.
47     iterations: Number of internal solver iterations.
48     blending_factor: Weight (0 to 1) for blending ORM+ output with the
49     smoothed best pure strategy. A factor of 0 means pure ORM+; 1 means
50     pure smoothed best pure strategy.
51     temperature: Temperature for softmax smoothing when calculating the
52     smoothed best pure strategy. Lower values make the smoothing sharper.
53     momentum_beta: Momentum parameter for optimistic updates to payoff gains.
54     gain_normalization: If True, normalizes payoff gains to make learning rate
55     more robust across games.
56     diversity_bonus_coeff: Coefficient for diversity bonus, encouraging
57     exploration of less-chosen policies.

```

```

58     return_average_strategy: If True, returns time-averaged strategies.
59     If False, returns last-iterate strategies.
60
61 Returns:
62     A list of mixed-strategies, one for each player, as numpy arrays.
63 """
64 num_players = len(meta_games)
65 num_strats = [m.shape[i] for i, m in enumerate(meta_games)]
66
67 if any(n_s == 0 for n_s in num_strats):
68     return [np.array([]).tolist() for _ in range(num_players)]
69
70 strategies = [np.ones(s, dtype=float) / s for s in num_strats]
71 cum_regrets = [np.zeros(s, dtype=float) for s in num_strats]
72 avg_strategies = [np.zeros(s, dtype=float) for s in num_strats]
73 prev_centered_payoff_gains = [np.zeros(s, dtype=float) for s in num_strats]
74
75 for t in range(iterations):
76     current_centered_payoff_gains = [np.zeros(s, dtype=float) for s in num_strats]
77     orm_strategies_this_iter = [np.zeros(s, dtype=float) for s in num_strats]
78
79     for p in range(num_players):
80         payoff_vec = meta_games[p]
81         for other_p in reversed(range(num_players)):
82             if other_p != p:
83                 payoff_vec = np.tensordot(payoff_vec, strategies[other_p], axes=([other_p], [0]))
84
85         centered_payoff_gains = payoff_vec - np.mean(payoff_vec)
86         current_centered_payoff_gains[p] = centered_payoff_gains
87
88         optimistic_payoff_gains = (1 + momentum_beta) * centered_payoff_gains - \
89             momentum_beta * prev_centered_payoff_gains[p]
90
91         diversity_bonus = diversity_bonus_coeff * (1.0 - strategies[p])
92
93         gains_for_regret_update = optimistic_payoff_gains + diversity_bonus
94
95         if gain_normalization:
96             max_abs_gain = np.max(np.abs(gains_for_regret_update))
97             if max_abs_gain > 1e-8:
98                 gains_for_regret_update /= max_abs_gain
99
100         cum_regrets[p] += gains_for_regret_update
101         cum_regrets[p] = np.maximum(0, cum_regrets[p])
102
103         sum_pos_regret = cum_regrets[p].sum()
104         if sum_pos_regret > 1e-12:
105             orm_strategies_this_iter[p] = cum_regrets[p] / sum_pos_regret
106         else:
107             orm_strategies_this_iter[p] = np.ones(num_strats[p]) / num_strats[p]
108
109         smoothed_best_pure = _smoothed_best_pure_strategy(payoff_vec, temperature)
110
111         strategies[p] = (1 - blending_factor) * orm_strategies_this_iter[p] + \
112             blending_factor * smoothed_best_pure
113
114         prev_centered_payoff_gains[p] = current_centered_payoff_gains[p]
115
116         if return_average_strategy: # Accumulate blended strategy only if average is requested
117             avg_strategies[p] += strategies[p]
118
119 if return_average_strategy:
120     final_strategies = []
121     for p in range(num_players):
122         sum_avg_strat = np.sum(avg_strategies[p])
123         if sum_avg_strat > 0:
124             final_strategies.append(avg_strategies[p] / sum_avg_strat)
125         else:
126             final_strategies.append(np.ones(num_strats[p]) / num_strats[p])
127     return final_strategies
128 else:
129     # If not returning average, return the last-iterate strategies
130     return strategies
131
132
133 class TrainMetaStrategySolver:
134     """A hybrid meta-solver for training that blends ORM+ with smoothed best pure strategies.
135

```

```

136 This solver aims to accelerate convergence to low-exploitable strategies by
137 dynamically balancing regret-minimization with a pull towards high-performing
138 (but smoothed) pure strategies. Optimistic updates, gain normalization, and
139 a diversity bonus are incorporated for improved learning dynamics. The
140 blending factor, temperature, and diversity bonus are annealed over the
141 outer PSRO iterations.
142 """
143
144 def __init__(self,
145             base_solver_iterations=1000, # Base number of internal iterations
146             iterations_per_policy_scale=20, # How much iterations scale per added policy
147             max_solver_iterations=5000, # Max internal solver iterations
148             initial_blending_factor=0.3, final_blending_factor=0.05,
149             initial_temperature=0.5, final_temperature=0.01,
150             momentum_beta=0.5,
151             gain_normalization=True,
152             initial_diversity_bonus_coeff=0.05,
153             final_diversity_bonus_coeff=0.001,
154             max_psro_iterations_for_annealing=75):
155     """Initializes hybrid ORM solver parameters for training.
156
157     Args:
158         base_solver_iterations: Base number of internal solver iterations for _hybrid_orm_solver.
159         iterations_per_policy_scale: Amount to increase internal solver iterations per added policy.
160         max_solver_iterations: Maximum internal solver iterations.
161         initial_blending_factor: Initial weight for the smoothed best pure strategy component.
162         final_blending_factor: Final weight for the smoothed best pure strategy component.
163         initial_temperature: Initial temperature for softmax smoothing.
164         final_temperature: Final temperature for softmax smoothing.
165         momentum_beta: Momentum for optimistic updates.
166         gain_normalization: Normalizes gains for scale-invariance.
167         initial_diversity_bonus_coeff: Max initial diversity bonus coefficient.
168         final_diversity_bonus_coeff: Min initial diversity bonus coefficient across PSRO iterations.
169         max_psro_iterations_for_annealing: PSRO iterations over which outer annealing occurs.
170     """
171     self._base_solver_iterations = base_solver_iterations
172     self._iterations_per_policy_scale = iterations_per_policy_scale
173     self._max_solver_iterations = max_solver_iterations
174     self._initial_blending_factor = initial_blending_factor
175     self._final_blending_factor = final_blending_factor
176     self._initial_temperature = initial_temperature
177     self._final_temperature = final_temperature
178     self._momentum_beta = momentum_beta
179     self._gain_normalization = gain_normalization
180     self._initial_diversity_bonus_coeff = initial_diversity_bonus_coeff
181     self._final_diversity_bonus_coeff = final_diversity_bonus_coeff
182     self._max_psro_iterations_for_annealing = max_psro_iterations_for_annealing
183     self._current_psro_iteration = 0
184
185 def get_meta_strategy(self, game, policy_sets, meta_games):
186     """Returns blended meta strategies for training.
187
188     Args:
189         game: The pypspiel game object.
190         policy_sets: A list of lists of policies, one list per player.
191             policy_sets[p][i] is player p's i-th policy. len(policy_sets[p]) ==
192             meta_games[0].shape[p].
193         meta_games: A list of n-dimensional numpy arrays, one per player. Each
194             array has shape (num_strats_p0, num_strats_p1, ..., num_strats_pn-1) and
195             meta_games[p][i0, i1, ..., in-1] is the payoff of player p when player k
196             chooses strategy ik.
197
198     Returns:
199         A list of blended mixed-strategies.
200     """
201     del game, policy_sets # Unused
202
203     self._current_psro_iteration += 1
204     current_psro_iter = self._current_psro_iteration
205
206     # Adaptive solver iterations: scale with current population size
207     num_current_policies_p0 = len(meta_games[0]) # Assuming symmetric populations
208     solver_iterations = int(self._base_solver_iterations +
209                             self._iterations_per_policy_scale * (num_current_policies_p0 - 1))
210     solver_iterations = np.clip(solver_iterations, self._base_solver_iterations, self._max_solver_iterations)
211
212     annealing_progress = min(1.0, current_psro_iter / self._max_psro_iterations_for_annealing)
213

```

```

214     blending_factor = (self._initial_blending_factor * (1.0 - annealing_progress) +
215                       self._final_blending_factor * annealing_progress)
216
217     temperature = (self._initial_temperature * (1.0 - annealing_progress) +
218                  self._final_temperature * annealing_progress)
219
220     diversity_bonus_coeff = (self._initial_diversity_bonus_coeff * (1.0 - annealing_progress) +
221                             self._final_diversity_bonus_coeff * annealing_progress)
222
223     blending_factor = np.clip(blending_factor, self._final_blending_factor, self._initial_blending_factor)
224     temperature = np.clip(temperature, self._final_temperature, self._initial_temperature)
225     diversity_bonus_coeff = np.clip(diversity_bonus_coeff, self._final_diversity_bonus_coeff, self._initial_diversity_bonus_coeff)
226
227     strategies = _hybrid_orm_solver(
228         meta_games,
229         iterations=solver_iterations, # Use adaptive iterations
230         blending_factor=blending_factor,
231         temperature=temperature,
232         momentum_beta=self._momentum_beta,
233         gain_normalization=self._gain_normalization,
234         diversity_bonus_coeff=diversity_bonus_coeff,
235         return_average_strategy=True # Training always uses averaged strategies for stability
236     )
237     return [s.tolist() for s in strategies]
238
239
240 class EvalMetaStrategySolver:
241     """Returns meta strategies for evaluation in PSRO.
242
243     This solver uses a hybrid approach, blending Optimistic Regret Matching+
244     with a smoothed best pure strategy, tailored for robust and accurate
245     exploitability measurement. The parameters are set to emphasize
246     exploitation for evaluation purposes, including optimistic updates and
247     gain normalization for stability, while keeping diversity bonus minimal.
248     Crucially, it returns the *last-iterate* strategy for a reactive estimate
249     of exploitability.
250     """
251
252     def __init__(self,
253                 base_solver_iterations=8000, # Base number of internal iterations
254                 iterations_per_policy_scale=50, # How much iterations scale per added policy
255                 max_solver_iterations=15000, # Max internal solver iterations
256                 blending_factor=0.01,
257                 temperature=0.001,
258                 momentum_beta=0.2,
259                 gain_normalization=True,
260                 diversity_bonus_coeff=0.0):
261         """Initializes hybrid ORM solver parameters for evaluation meta-strategies.
262
263         Args:
264             base_solver_iterations: Base number of internal solver iterations for _hybrid_orm_solver.
265             iterations_per_policy_scale: Amount to increase internal solver iterations per added policy.
266             max_solver_iterations: Maximum internal solver iterations.
267             blending_factor: Weight (0 to 1) for the smoothed best pure strategy component.
268             temperature: Temperature for softmax smoothing.
269             momentum_beta: Momentum for optimistic updates.
270             gain_normalization: Normalizes gains for scale-invariance.
271             diversity_bonus_coeff: Diversity bonus, kept very low for evaluation.
272         """
273         self._base_solver_iterations = base_solver_iterations
274         self._iterations_per_policy_scale = iterations_per_policy_scale
275         self._max_solver_iterations = max_solver_iterations
276         self._blending_factor = blending_factor
277         self._temperature = temperature
278         self._momentum_beta = momentum_beta
279         self._gain_normalization = gain_normalization
280         self._diversity_bonus_coeff = diversity_bonus_coeff
281
282     def get_meta_strategy(self, game, policy_sets, meta_games):
283         """Returns blended meta strategies for evaluation in policy-space response oracles.
284
285         Args:
286             game: The pypiel game object.
287             policy_sets: A list of lists of policies, one list per player.
288                 policy_sets[p][i] is player p's i-th policy. len(policy_sets[p]) ==
289                 meta_games[0].shape[p].
290             meta_games: A list of n-dimensional numpy arrays, one per player. Each
291                 array has shape (num_strats_p0, num_strats_p1, ..., num_strats_pn-1) and

```

```

292     meta_games[p][i0, i1, ..., in-1] is the payoff of player p when player k
293     chooses strategy ik.
294
295 Returns:
296     A list of mixed-strategies, one for each player. Each mixed strategy is
297     a list of non-negative weights (not necessarily normalized). It is used
298     for evaluation of the current PSRO policies. E.g., computing
299     exploitability.
300 """
301 del game, policy_sets # Unused
302
303 num_current_policies_p0 = len(meta_games[0]) # Assuming symmetric populations
304 solver_iterations = int(self._base_solver_iterations +
305                        self._iterations_per_policy_scale * (num_current_policies_p0 - 1))
306 solver_iterations = np.clip(solver_iterations, self._base_solver_iterations, self._max_solver_iterations)
307
308 strategies = _hybrid_orm_solver(
309     meta_games,
310     iterations=solver_iterations, # Use adaptive iterations
311     blending_factor=self._blending_factor,
312     temperature=self._temperature,
313     momentum_beta=self._momentum_beta,
314     gain_normalization=self._gain_normalization,
315     diversity_bonus_coeff=self._diversity_bonus_coeff,
316     return_average_strategy=False # Eval explicitly requests last-iterate strategy
317 )
318 return [s.tolist() for s in strategies]

```

## 7.2 Asymmetric Optimistic Discounted CFR

In an early trial, we discover another variant of CFR by training on 2-player Kuhn Poker, 2-player Leduc Poker, 4-Card Goofspiel and 4-side Liar Dice. It also shows good empirical performance and employs relatively more conventional mechanisms than VAD-CFR. We show its empirical performance together with other CFR variants in Figure 3. The source code is in Listing 7. Unlike manual heuristics which often rely on symmetric or fixed update rules, AOD-CFR utilizes a set of dynamic, asymmetric mechanisms discovered via evolutionary search.

- **Adaptive Regret Discounting:** The algorithm employs a linear schedule for discounting cumulative regrets. The positive regret discount factor,  $\alpha$ , transitions from 1.0  $\rightarrow$  2.5 over 500 iterations, while the negative regret discount factor,  $\beta$ , transitions from 0.5  $\rightarrow$  0.0. This allows the algorithm to aggressively prune suboptimal actions early (via  $\beta$ ) while increasingly emphasizing recent information for optimal actions (via  $\alpha$ ).
- **Asymmetric Instantaneous Scaling:** A distinct feature of AOD-CFR is its sign-dependent scaling of instantaneous regret. As revealed in the source code, the update rule scales instantaneous regret by 1.1 when both cumulative and instantaneous regrets are positive, but attenuates it by 0.9 when cumulative regret is positive and instantaneous regret is negative. This asymmetry likely functions as a momentum-preservation mechanism, reinforcing established beneficial actions while damping noise.
- **Trend-Based Policy Optimism:** In the policy derivation step, AOD-CFR incorporates a “trend-based optimism” term. It tracks an Exponential Moving Average (EMA) of cumulative regrets (decay  $\approx$  0.1) and adds a scaled deviation term to the accumulated regret before applying Regret Matching. Specifically, the code implements a scaling factor `optimism_trend_scale` applied to the difference between current cumulative regret and its EMA.
- **Polynomial Policy Averaging:** The average strategy is computed using an adaptive polynomial weighting scheme where the exponent  $\gamma$  scales linearly from 1.0  $\rightarrow$  5.0.

### Listing 7: AOD-CFR

```

1 import math
2 import collections
3 import attr
4
5 class RegretAccumulator:
6     """Updates cumulative regret with adaptive positive and negative discounting,
7     and asymmetric instantaneous regret scaling."""
8
9     def __init__(self,
10                 alpha_start: float = 1.0, alpha_max: float = 2.5, schedule_T_alpha: float = 500.0,
11                 beta_start: float = 0.5, beta_max: float = 0.0, schedule_T_beta: float = 500.0,
12                 pos_cum_pos_inst_scale: float = 1.1, # R_cum > 0, r_inst > 0
13                 pos_cum_neg_inst_scale: float = 0.9, # R_cum > 0, r_inst < 0
14                 neg_cum_pos_inst_scale: float = 0.7, # R_cum < 0, r_inst > 0
15                 neg_cum_neg_inst_scale: float = 1.2): # R_cum < 0, r_inst < 0
16         """Initializes the regret accumulator.
17
18     Args:
19         alpha_start: Initial exponent for positive cumulative regret discounting.
20         alpha_max: Final exponent for positive cumulative regret discounting.

```

```

21     schedule_T_alpha: Iterations for alpha transition.
22     beta_start: Initial exponent for negative cumulative regret discounting.
23     beta_max: Final exponent for negative cumulative regret discounting.
24     schedule_T_beta: Iterations for beta transition.
25     pos_cum_pos_inst_scale: Scaling factor for instantaneous regret when
26         cumulative regret is positive and instantaneous is positive.
27     pos_cum_neg_inst_scale: Scaling factor for instantaneous regret when
28         cumulative regret is positive and instantaneous is negative.
29     neg_cum_pos_inst_scale: Scaling factor for instantaneous regret when
30         cumulative regret is negative and instantaneous is positive.
31     neg_cum_neg_inst_scale: Scaling factor for instantaneous regret when
32         cumulative regret is negative and instantaneous is negative.
33     """
34     if not (alpha_start >= 0 and alpha_max >= 0 and schedule_T_alpha >= 0):
35         raise ValueError("Alpha parameters and schedule must be non-negative.")
36     if not (beta_start >= 0 and beta_max >= 0 and schedule_T_beta >= 0):
37         raise ValueError("Beta parameters and schedule must be non-negative.")
38     if not (pos_cum_pos_inst_scale >= 0 and pos_cum_neg_inst_scale >= 0 and
39             neg_cum_pos_inst_scale >= 0 and neg_cum_neg_inst_scale >= 0):
40         raise ValueError("Instantaneous regret scaling factors must be non-negative.")
41
42     self._alpha_start = alpha_start
43     self._alpha_max = alpha_max
44     self._schedule_T_alpha = max(1.0, schedule_T_alpha)
45     self._beta_start = beta_start
46     self._beta_max = beta_max
47     self._schedule_T_beta = max(1.0, schedule_T_beta)
48
49     self._pos_cum_pos_inst_scale = pos_cum_pos_inst_scale
50     self._pos_cum_neg_inst_scale = pos_cum_neg_inst_scale
51     self._neg_cum_pos_inst_scale = neg_cum_pos_inst_scale
52     self._neg_cum_neg_inst_scale = neg_cum_neg_inst_scale
53     self._epsilon = 1e-9 # Small value for sign comparisons
54
55
56 def update_accumulate_regret(
57     self, info_state_node, iteration_number, cfr_regrets
58 ):
59     """Updates cumulative regrets with adaptive discounting and asymmetric instantaneous regret scaling.
60
61     Args:
62         info_state_node: Data structure for the information set.
63         iteration_number: The current CFR iteration (0-based).
64         cfr_regrets: Instantaneous counterfactual regrets for the current iteration.
65     Returns:
66         Updated cumulative regret dictionary for each action.
67     """
68     t = float(iteration_number) + 1.0 # 1-indexed for weighting calculations
69     iter_float = float(iteration_number) # 0-based for schedule progress
70
71     # Calculate current_alpha for positive cumulative regret discounting
72     t_norm_alpha = min(1.0, iter_float / self._schedule_T_alpha)
73     current_alpha = self._alpha_start + (self._alpha_max - self._alpha_start) * t_norm_alpha
74
75     # Calculate current_beta for negative cumulative regret discounting
76     t_norm_beta = min(1.0, iter_float / self._schedule_T_beta)
77     current_beta = self._beta_start + (self._beta_max - self._beta_start) * t_norm_beta
78
79     updated_cumulative_regret = {}
80     legal_actions = info_state_node.legal_actions
81
82     for action in legal_actions:
83         prev_cumulative_regret = info_state_node.cumulative_regret.get(action, 0.0)
84         current_instantaneous_regret = cfr_regrets.get(action, 0.0)
85
86         # Determine instantaneous regret scaling factor based on signs
87         instantaneous_scaling_factor = 1.0 # Default
88
89         if prev_cumulative_regret > self._epsilon: # Cumulative regret is positive
90             if current_instantaneous_regret > self._epsilon:
91                 instantaneous_scaling_factor = self._pos_cum_pos_inst_scale
92             elif current_instantaneous_regret < -self._epsilon:
93                 instantaneous_scaling_factor = self._pos_cum_neg_inst_scale
94         elif prev_cumulative_regret < -self._epsilon: # Cumulative regret is negative
95             if current_instantaneous_regret > self._epsilon:
96                 instantaneous_scaling_factor = self._neg_cum_pos_inst_scale
97             elif current_instantaneous_regret < -self._epsilon:
98                 instantaneous_scaling_factor = self._neg_cum_neg_inst_scale

```

```

99     # If one or both are near zero, instantaneous_scaling_factor remains 1.0
100
101     scaled_instantaneous_regret = current_instantaneous_regret * instantaneous_scaling_factor
102
103     # Calculate discount factor for the previous cumulative regret based on its sign
104     discount_factor = 0.0
105     if prev_cumulative_regret > 0:
106         try:
107             t_pow_alpha = math.pow(t, current_alpha)
108             # If t^alpha overflows, the factor approaches 1.0
109             discount_factor = t_pow_alpha / (t_pow_alpha + 1.0) if t_pow_alpha != float('inf') else 1.0
110         except ValueError: # Fallback for potential issues with very large t and alpha
111             discount_factor = 1.0 if t > 1.1 else 0.5 # t=1 is 0.5 (1/2), large t is 1.0 (inf/inf)
112     else: # prev_cumulative_regret is non-positive
113         try:
114             t_pow_beta = math.pow(t, current_beta)
115             # If t^beta overflows, the factor approaches 1.0
116             discount_factor = t_pow_beta / (t_pow_beta + 1.0) if t_pow_beta != float('inf') else 1.0
117         except ValueError: # Fallback for potential issues with very large t and beta
118             discount_factor = 1.0 if t > 1.1 else 0.5 # t=1 is 0.5 (1/2), large t is 1.0 (inf/inf)
119
120     discounted_prev_regret = prev_cumulative_regret * discount_factor
121
122     # Add the scaled current instantaneous regret.
123     new_cumulative_regret = discounted_prev_regret + scaled_instantaneous_regret
124
125     updated_cumulative_regret[action] = new_cumulative_regret
126
127     return updated_cumulative_regret
128
129
130 import math
131 import collections
132 import attr # Ensure attr is imported if not already globally
133
134
135 class PolicyFromRegretAccumulator:
136     """Calculates the current policy using cumulative regrets plus an optimistic term
137     based on the deviation of current cumulative regret from its EMA."""
138
139     def __init__(self,
140                 use_squared_weights: bool = True,
141                 cumulative_regret_ema_alpha: float = 0.1, # EMA factor for cumulative regrets (lower means more smoothing)
142                 optimism_trend_scale: float = 0.5): # Scaling factor for the R^t - EMA(R^{t-1}) term
143         """Initializes the policy calculator with trend-based optimism.
144
145         Args:
146             use_squared_weights: If True (default), use max(0, R + scaled_optimism)^2 for policy weights (CFR+ style).
147                                 If False, use max(0, R + scaled_optimism) for policy weights (RM+ style).
148             cumulative_regret_ema_alpha: Decay factor for the EMA of cumulative regrets (0 < alpha <= 1).
149                                         Higher alpha gives more weight to recent cumulative regrets.
150             optimism_trend_scale: Scaling factor applied to the difference between current cumulative regret
151                                 and its EMA (R^t - EMA(R^{t-1})) before adding it to R^t.
152
153         """
154         if not (0.0 < cumulative_regret_ema_alpha <= 1.0):
155             raise ValueError("cumulative_regret_ema_alpha must be between 0 (exclusive) and 1 (inclusive).")
156         if optimism_trend_scale < 0:
157             raise ValueError("optimism_trend_scale must be non-negative.")
158
159         self._use_squared_weights = use_squared_weights
160         self._cumulative_regret_ema_alpha = cumulative_regret_ema_alpha
161         self._optimism_trend_scale = optimism_trend_scale
162         # Store EMA of *cumulative* regrets: {info: set_index: {action: ema_value}}
163         self._ema_cumulative_regrets = collections.defaultdict(lambda: collections.defaultdict(float))
164         # Epsilon for safe floating point comparisons for policy normalization.
165         self._epsilon_norm = 1e-12
166
167     def get_updated_current_policy(self, info_state_node, iteration_number, cfr_regrets, previous_policy):
168         """Calculates the current policy using cumulative regrets plus trend-based optimism.
169
170         Args:
171             info_state_node: a data structure corresponding to an information set with
172                             cumulative_regret and cumulative_policy stored.
173                             cumulative_regret (R^t) has been updated by RegretAccumulator
174                             for the current iteration.
175             iteration_number: the current CFR iteration (0-based) (unused but kept for API).
176             cfr_regrets: the instantaneous counterfactual regrets (r^t) for the current iteration (unused).

```

```

177     previous_policy: the previous policy at the current information set (unused).
178 Returns:
179     updated current policy dictionary at the current information set.
180 """
181 cumulative_regrets_map = info_state_node.cumulative_regret # R^t
182 actions = info_state_node.legal_actions
183 infoset_index = info_state_node.index_in_tabular_policy # Unique key for EMA storage
184
185 action_to_policy_weight = {}
186 sum_policy_weights = 0.0
187 action_to_new_ema = {} # Store new EMA values before overwriting state
188
189 if not actions: # Handle cases with no legal actions
190     return {}
191
192 num_actions = len(actions) # Used for uniform policy fallback
193
194 # Get the EMA map for this specific infoset
195 infoset_ema_map = self._ema_cumulative_regrets[infoset_index]
196
197 for action in actions:
198     # Get the current cumulative regret R^t(a)
199     current_cum_regret = cumulative_regrets_map.get(action, 0.0)
200     # Get the previous EMA of cumulative regret EMA(R^{t-1})(a)
201     prev_ema_cum_regret = infoset_ema_map.get(action, 0.0) # Defaults to 0.0 on first encounter
202
203     # Calculate the deviation = R^t(a) - EMA(R^{t-1})(a)
204     regret_trend = current_cum_regret - prev_ema_cum_regret
205
206     # Calculate the scaled optimistic term
207     scaled_optimism = self._optimism_trend_scale * regret_trend
208
209     # Calculate the value used for Regret Matching: R^t(a) + scaled_optimism
210     combined_value = current_cum_regret + scaled_optimism
211
212     # Apply Regret Matching logic: take the positive part.
213     positive_combined_value = max(0.0, combined_value)
214
215     # Use linear or squared weighting based on the flag
216     if self._use_squared_weights:
217         policy_weight = positive_combined_value ** 2 # CFR+ style
218     else:
219         policy_weight = positive_combined_value # RM+ style
220
221     action_to_policy_weight[action] = policy_weight
222     sum_policy_weights += policy_weight
223
224     # Update the EMA of cumulative regret for the next iteration (t)
225     # EMA^t = alpha * R^t + (1 - alpha) * EMA^{t-1}
226     new_ema_cum_regret = self._cumulative_regret_ema_alpha * current_cum_regret + \
227         (1.0 - self._cumulative_regret_ema_alpha) * prev_ema_cum_regret
228     action_to_new_ema[action] = new_ema_cum_regret
229
230
231 # Update the stored EMA values *after* using the previous values for all actions
232 for action in actions:
233     infoset_ema_map[action] = action_to_new_ema[action]
234
235 # Normalize to get the policy
236 info_state_policy = {}
237 if sum_policy_weights > self._epsilon_norm: # Use normalization epsilon
238     # Normalize positive weights (or their squares) to get policy probabilities
239     for action in actions:
240         info_state_policy[action] = action_to_policy_weight[action] / sum_policy_weights
241 else:
242     # Default to uniform policy if sum of weights is zero or numerically close to zero
243     uniform_prob = 1.0 / num_actions if num_actions > 0 else 0.0
244     for action in actions:
245         info_state_policy[action] = uniform_prob
246
247 return info_state_policy
248
249
250 class PolicyAccumulator:
251     """Updates cumulative policy using adaptive polynomial weighting (like DCFR/PDCFR+)."""
252
253     def __init__(self,
254                 gamma_start: float = 1.0, gamma_max: float = 5.0, gamma_schedule_T: float = 500.0):

```

```

255     """Initializes the policy accumulator with adaptive gamma.
256
257     Args:
258         gamma_start: Initial exponent for weighting the current policy (t^gamma).
259         gamma_max: Max gamma value for later exploitation.
260         gamma_schedule_T: Number of iterations over which gamma transitions.
261     """
262     if gamma_start < 0 or gamma_max < 0 or gamma_schedule_T < 0:
263         raise ValueError("PolicyAccumulator gamma parameters must be non-negative.")
264     self._gamma_start = gamma_start
265     self._gamma_max = gamma_max
266     self._gamma_schedule_T = max(1.0, gamma_schedule_T) # Avoid division by zero
267
268
269 def update_accumulate_policy(
270     self,
271     info_state_node,
272     iteration_number,
273     info_state_policy,
274     cfr_regrets, # Not used
275     reach_prob,
276     counterfactual_reach_prob, # Not used
277 ):
278     """Updates the cumulative policy using adaptive polynomial weighting.
279
280     Args:
281         info_state_node: Data structure for the information set.
282         iteration_number: The current CFR iteration (0-based).
283         info_state_policy: The current policy profile for this iteration.
284         reach_prob: Player's reach probability for this iteration.
285     Returns:
286         Updated cumulative policy dictionary for each action.
287     """
288     # Use 1-based iteration number t for calculations
289     t = float(iteration_number) + 1.0
290     # Use 0-based iteration for schedule progress calculation
291     iter_float = float(iteration_number)
292
293     # Calculate current_gamma based on linear schedule
294     t_norm_gamma = min(1.0, iter_float / self._gamma_schedule_T)
295     current_gamma = self._gamma_start + (self._gamma_max - self._gamma_start) * t_norm_gamma
296
297     # Calculate t^current_gamma, which is used for both discounting and weighting.
298     t_pow_gamma = 1.0
299     if t > 0: # t is 1-indexed (float(iteration_number) + 1.0), so t=0 is not possible here.
300         try:
301             t_pow_gamma = math.pow(t, current_gamma)
302         except ValueError: # Fallback for extremely large exponents to avoid overflow, treat as very large.
303             t_pow_gamma = float('inf')
304
305     # Calculate the discount factor for the previous cumulative policy sum.
306     # This mirrors the discounting logic used for positive cumulative regrets in RegretAccumulator.
307     # If t_pow_gamma is 'inf', discount factor approaches 1.0.
308     policy_discount_factor = t_pow_gamma / (t_pow_gamma + 1.0) if t_pow_gamma != float('inf') else 1.0
309
310     # Calculate the weight for the current policy's contribution (Adaptive CFR style).
311     # This weights the current policy by t^gamma and the player's reach probability.
312     current_policy_contribution_weight = t_pow_gamma * reach_prob
313
314     updated_cumulative_policy = {}
315     # Iterate over all legal actions and actions present in the cumulative policy
316     # to ensure proper discounting and updating of all relevant entries.
317     all_actions = set(info_state_node.legal_actions) | set(info_state_node.cumulative_policy.keys())
318
319     for action in all_actions:
320         previous_cumulative_policy = info_state_node.cumulative_policy.get(action, 0.0)
321         current_policy_prob = info_state_policy.get(action, 0.0)
322
323         # Discount the previously accumulated policy sum.
324         discounted_previous_cumulative_policy = previous_cumulative_policy * policy_discount_factor
325
326         # Add the current policy's weighted contribution.
327         current_contribution = current_policy_contribution_weight * current_policy_prob
328
329         updated_cumulative_policy[action] = (
330             discounted_previous_cumulative_policy + current_contribution
331         )
332

```

```

333 # Ensure all *legal* actions are present in the final dict, even if value is 0.0
334 for action in info_state_node.legal_actions:
335     if action not in updated_cumulative_policy:
336         # This happens if an action was legal but never accumulated policy before
337         # and isn't in current policy (prob=0). Its value should be 0
338         # after discounting the previous 0 and adding 0 contribution.
339         updated_cumulative_policy[action] = 0.0
340
341 return updated_cumulative_policy

```

### 7.3 Results on 11 Games

All Games Performance

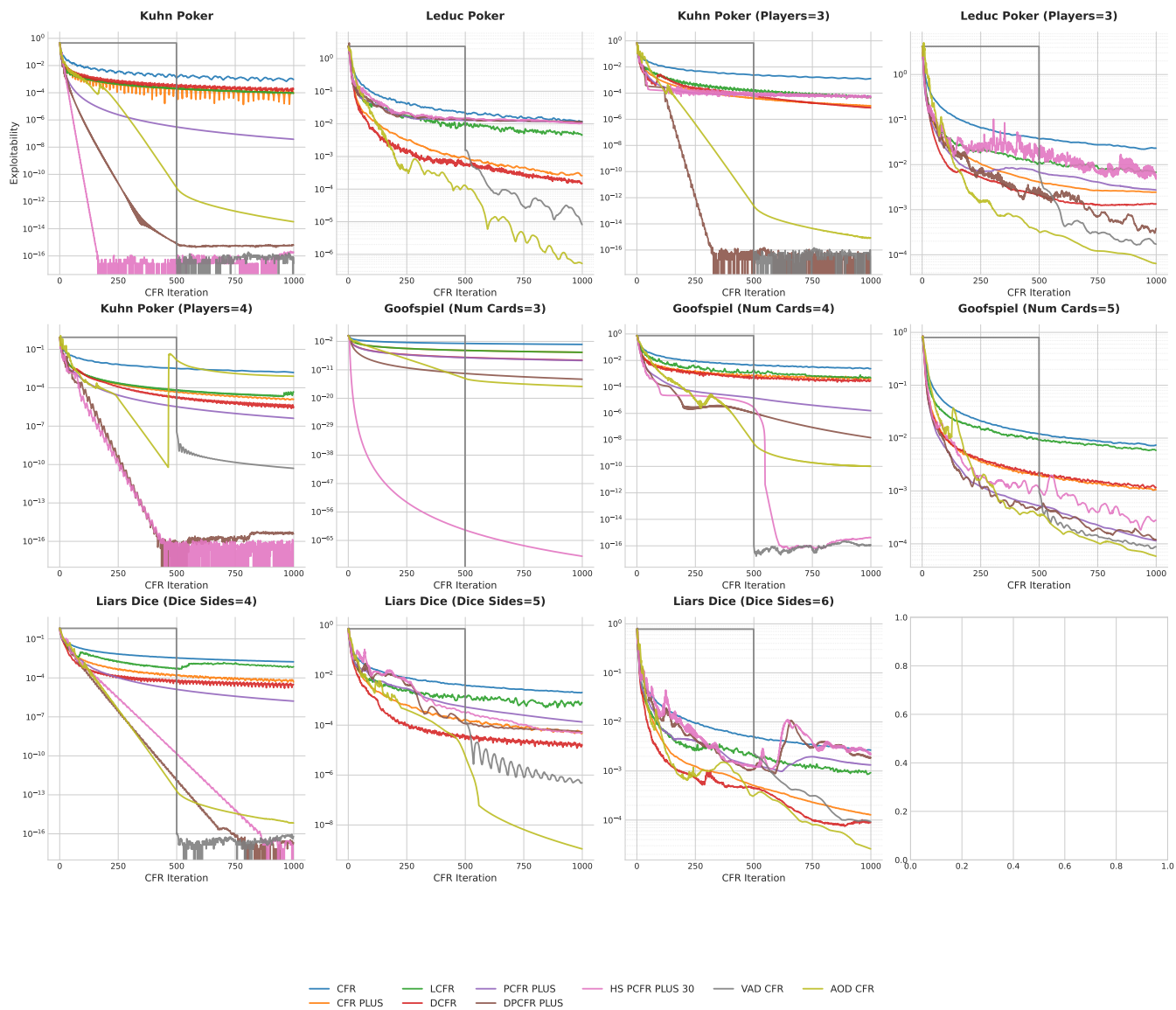


Figure 3: CFR variants performances on All Games.

## All Games Performance

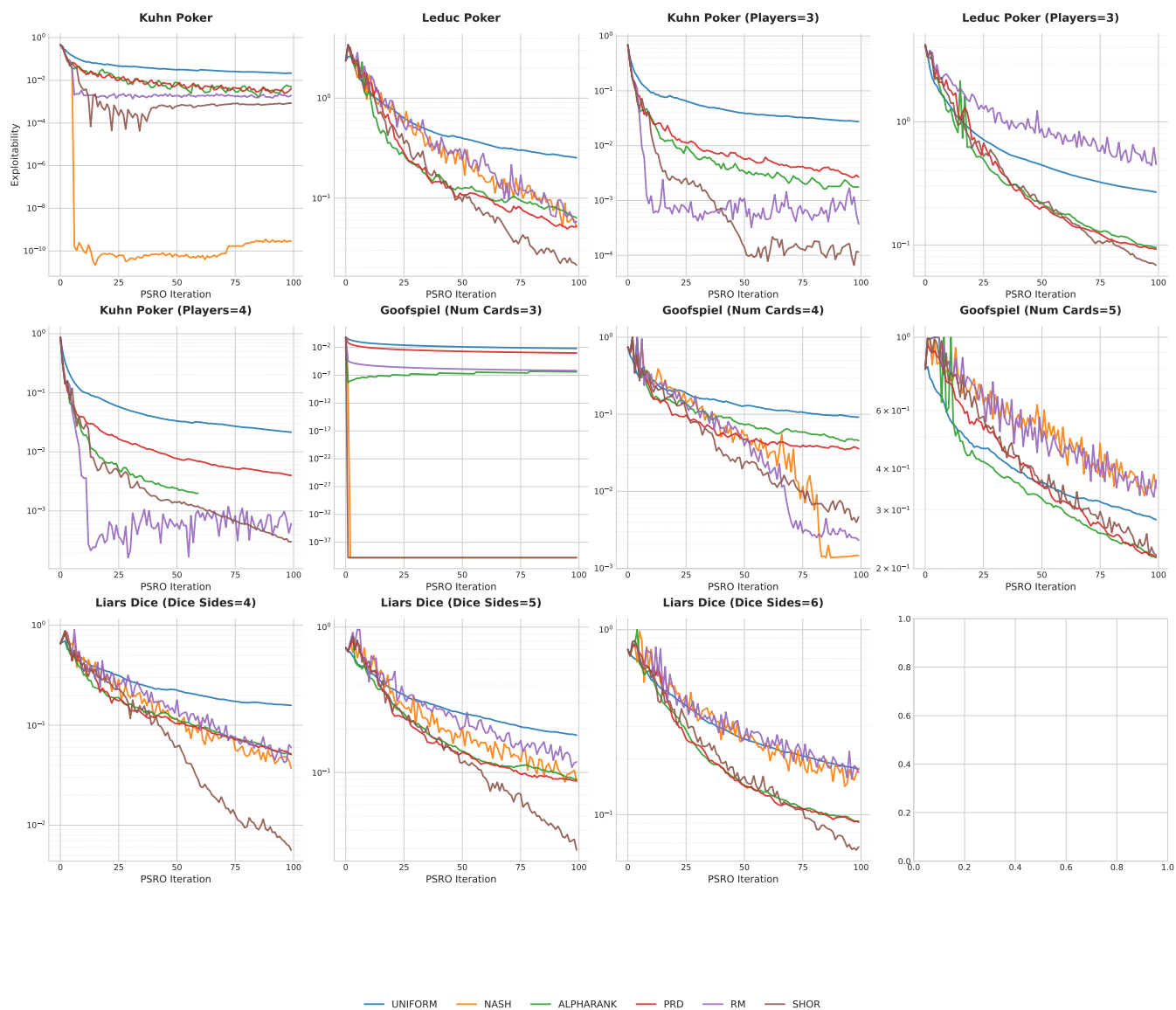


Figure 4: PSRO variants performances on All Games.

## 7.4 Prompts

### Listing 8: Prompt for Evolving CFR

```

1 Act as an expert in game theory, multiagent learning, online learning and optimization. Your task is to iteratively improve a new variant of
  counterfactual regret minimization. The primary goal is to speed up convergence to low-exploitable strategies.
2
3 Always adhere to best practices in Python coding.
4
5
6 A key data structure that is used is infostate_node:
7
8 ```python

```

```

9 @attr.s
10 class _InfoStateNode(object):
11     """An object wrapping values associated to an information state."""
12     # The list of the legal actions.
13     legal_actions = attr.ib()
14     index_in_tabular_policy = attr.ib()
15     # Map from information states string representations and actions to the
16     # counterfactual regrets, accumulated over the policy iterations
17     cumulative_regret = attr.ib(factory=lambda: collections.defaultdict(float))
18     # Same as above for the cumulative of the policy probabilities computed
19     # during the policy iterations
20     cumulative_policy = attr.ib(factory=lambda: collections.defaultdict(float))
21     ...
22
23 You are allowed to modify three key components of CFR: (1) how are the regret values accumulated at each info set (RegretAccumulator) (2) how
    to obtain a current policy at the current iteration from the current cumulative_regret (PolicyFromRegretAccumulator) and (3) how to
    accumulate policies across iterations to compute an average policy (PolicyAccumulator).
24
25 # Prior programs
26 Previously we found that the following programs performed well on the task at hand:
27
28 {previous_programs}
29
30 # Current program
31 Here is the current program we are trying to improve (you will need to propose a modification to it below):
32
33 {code}
34
35 # *SEARCH/REPLACE block* Rules:
36
37 Every *SEARCH/REPLACE block* must use this format:
38 1. The opening fence: ```python
39 2. The start of search block: <<<<<< SEARCH
40 3. A contiguous chunk of up to 4 lines to search for in the existing source code
41 4. The dividing line: =====
42 5. The lines to replace into the source code
43 6. The end of the replace block: >>>>>> REPLACE
44 7. The closing fence: ```
45
46 Every *SEARCH* section must *EXACTLY MATCH* the existing file content, character for character, including all comments, docstrings, etc.
47
48 *SEARCH/REPLACE* blocks will replace *all* matching occurrences.
49 Include enough lines to make the SEARCH blocks uniquely match the lines to change.
50
51 Keep *SEARCH/REPLACE* blocks concise.
52 Break large *SEARCH/REPLACE* blocks into a series of smaller blocks that each change a small portion of the file.
53 Include just the changing lines, and a few surrounding lines if needed for uniqueness.
54 Do not include long runs of unchanging lines in *SEARCH/REPLACE* blocks.
55
56 To move code within a file, use 2 *SEARCH/REPLACE* blocks: 1 to delete it from its current location, 1 to insert it in the new location.
57
58 Make sure that the changes you propose are consistent with each other. For example, if you refer to a new config variable somewhere, you
    should also propose a change to add that variable.
59
60 Example:
61 ```python
62 <<<<<< SEARCH
63     return total_loss
64 =====
65     # Add sparsity-promoting regularization to the loss.
66     total_loss += self.hypers.l1_reg_weight * l1_reg
67
68     return total_loss
69 {replace}
70 ```
71 and
72 ```python
73 <<<<<< SEARCH
74     return hyper.zipit([
75 =====
76     return hyper.zipit([
77         hyper.uniform('l1_reg_weight', hyper.interval(0.0, 0.01)),
78 {replace}
79 ```
80
81 {lazy_prompt}
82 ONLY EVER RETURN CODE IN A *SEARCH/REPLACE BLOCK*!
83

```

```

84 # Task
85 {task_instruction} {focus_sentence} {trigger_chain_of_thought}
86 Describe each change with a *SEARCH/REPLACE block*.

```

### Listing 9: Prompt for Evolving PSRO

```

1 Act as an expert in game theory, multiagent learning, online learning and optimization. Your task is to iteratively improve a variant of
  Policy-Space Response Oracles (PSRO). The primary goal is to speed up convergence to low-exploitable strategies.
2
3 # PSRO Overview
4
5 PSRO iteratively builds a population of policies for each player and computes meta-strategies (distributions over policies) to guide
  training and evaluation.
6
7 **Each iteration:**
8 1. **Empirical game**: Simulate payoffs between all policy combinations to form a game tensor.
9 2. **Train-time meta-strategy**: Compute a distribution over current policies for each player. This determines what opponents the best-
  response oracle trains against.
10 3. **Best response**: Add a new policy for each player that best responds to opponents' train-time meta-strategies.
11 4. **Eval-time meta-strategy**: Compute a (possibly different) distribution for evaluation, e.g., to measure exploitability.
12
13 **Your task**: Improve both the **train-time** and **eval-time** meta-strategy solvers. These serve different purposes: train-time guides
  population growth, eval-time assesses solution quality.
14
15 # Available Utilities
16
17 **Best Response**
18 - `BestResponsePolicy(game, player_id, policy)`: Returns a best-response policy for `player_id` against `policy`.
19 - Use `.value(game.new_initial_state())` to get the BR value.
20
21 **Policy Aggregation**
22 - `PolicyAggregator(game).aggregate(pids, policy_sets, weights)`: Creates a mixed policy from weighted pure policies.
23 - `pids`: `list(range(game.num_players()))`
24 - `weights`: `[weights_p0, weights_p1, ...]`, each matching `len(policy_sets[p])`
25
26 **Policy Evaluation**
27 - `expected_game_score.policy_value(state, joint_policy)`: Returns expected payoffs (list, one per player).
28 - `joint_policy`: `[policy_p0, policy_p1, ...]`
29
30 Always adhere to best practices in Python coding.
31
32 # Prior programs
33 Previously we found that the following programs performed well on the task at hand:
34
35 {previous_programs}
36
37 # Current program
38 Here is the current program we are trying to improve (you will need to propose a modification to it below):
39
40 {code}
41
42 # *SEARCH/REPLACE block* Rules:
43
44 Every *SEARCH/REPLACE block* must use this format:
45 1. The opening fence: ```python
46 2. The start of search block: <<<<<< SEARCH
47 3. A contiguous chunk of up to 4 lines to search for in the existing source code
48 4. The dividing line: =====
49 5. The lines to replace into the source code
50 6. The end of the replace block: >>>>>> REPLACE
51 7. The closing fence: ```
52
53 Every *SEARCH* section must *EXACTLY MATCH* the existing file content, character for character, including all comments, docstrings, etc.
54
55 *SEARCH/REPLACE* blocks will replace *all* matching occurrences.
56 Include enough lines to make the SEARCH blocks uniquely match the lines to change.
57
58 Keep *SEARCH/REPLACE* blocks concise.
59 Break large *SEARCH/REPLACE* blocks into a series of smaller blocks that each change a small portion of the file.
60 Include just the changing lines, and a few surrounding lines if needed for uniqueness.
61 Do not include long runs of unchanging lines in *SEARCH/REPLACE* blocks.
62
63 To move code within a file, use 2 *SEARCH/REPLACE* blocks: 1 to delete it from its current location, 1 to insert it in the new location.
64
65 Make sure that the changes you propose are consistent with each other. For example, if you refer to a new config variable somewhere, you
  should also propose a change to add that variable.
66
67 Example:

```

```
68 ```python
69 <<<<<< SEARCH
70     return total_loss
71 =====
72     # Add sparsity-promoting regularization to the loss.
73     total_loss += self.hypers.l1_reg_weight * l1_reg
74
75     return total_loss
76 {replace}
77 ```
78 and
79 ```python
80 <<<<<< SEARCH
81     return hyper.zipit([
82 =====
83     return hyper.zipit([
84         hyper.uniform('l1_reg_weight', hyper.interval(0.0, 0.01)),
85 {replace}
86 ```
87
88 {lazy_prompt}
89 ONLY EVER RETURN CODE IN A *SEARCH/REPLACE BLOCK*!
90
91 # Task
92 {task_instruction} {focus_sentence} {difficulty_hint} {exploration_nudge} {trigger_chain_of_thought}
93 Describe each change with a *SEARCH/REPLACE block*.
```