# Efficient Policy Space Response Oracles

**Ming Zhou**[*1,4]  **Jingxiao Chen**[1]  **Ying Wen**[1]  **Weinan Zhang**[1]

**Yaodong Yang**[2]  **Yong Yu**[1]  **Jun Wang**[3,4]

[1]Shanghai Jiao Tong University, [2]Institute for Artificial Intelligence, Peking University
[3]University College London,[4]Shanghai Digital Brain Laboratory

## Abstract

Policy Space Response Oracle methods (PSRO) provide a general solution to learn Nash equilibrium in two-player zero-sum games but suffer from two drawbacks: (1) the *computation inefficiency* due to the need for consistent meta-game evaluation via simulations, and (2) the *exploration inefficiency* due to finding the best response against a fixed meta-strategy at every epoch. In this work, we propose Efficient PSRO (EPSRO) that largely improves the efficiency of the above two steps. Central to our development is the newly-introduced subroutine of *no-regret optimization* on the *unrestricted-restricted (URR)* game. By solving URR at each epoch, one can evaluate the current game and compute the best response in one forward pass without the need for meta-game simulations. Theoretically, we prove that the solution procedures of EPSRO offer a monotonic improvement on the exploitability, which none of existing PSRO methods possess. Furthermore, we prove that the no-regret optimization has a regret bound of $\mathcal{O}(\sqrt{T \log{[(k^2 + k)/2]}})$, where $k$ is the size of restricted policy set. Most importantly, a desirable property of EPSRO is that it is parallelizable, this allows for highly efficient exploration in the policy space that induces behavioral diversity. We test EPSRO on three classes of games, and report a 50x speedup in wall-time and 10x data efficiency while maintaining similar exploitability as existing PSRO methods on Kuhn and Leduc Poker games.

## 1 Introduction

Policy Space Response Oracles (PSRO) [19] is a general multi-agent reinforcement learning algorithm, which has been applied in many non-trivial multi-agent learning tasks [35, 3, 23]. In general, PSRO aims to find an approximate Nash equilibrium (NE) by iteratively expanding a restricted game formed by a set of restricted policy sets, which is ideally much smaller than the original game. At each epoch, PSRO executes sequential learning composed of a meta-game solving and a learning of best responses. Though PSRO does not need to learn policies in the original game directly, the learning of PSRO is still inefficient in solving meta-game and learning high-quality best responses.

Specifically, PSRO is **computation inefficient** to solve a geometrically growing meta-game because it relies on numerous simulations across the Cartesian space of growing policy sets [28, 39]. Moreover, learning against a fixed opponent meta-strategy to find a best response is **exploration inefficient**. In such a way, PSRO has no non-degenerate guarantee on the expansion of restricted policy sets, since the fixed meta-strategy is only a best response to the restricted policy set [36]. Despite playing against fixed opponent meta-strategies can theoretically expand the policy sets [24], it has no guarantee that the discovered best responses can still hold the strength when opponents deviate their strategies.

---

[*]Corresponding author: mingak@sjtu.edu.cn

Therefore, PSRO needs to add all possible policies from the original game and generates large restricted policy sets in the worst case, making it slow to converge [24]. A straightforward idea to improve the learning efficiency is to utilize parallelism for the learning of best responses [19, 2, 24, 22]. However, all of these existing methods still require simulations to solve meta-game and learn best responses by playing against fixed meta-strategies.



(a) Player 1's URR game dynamics

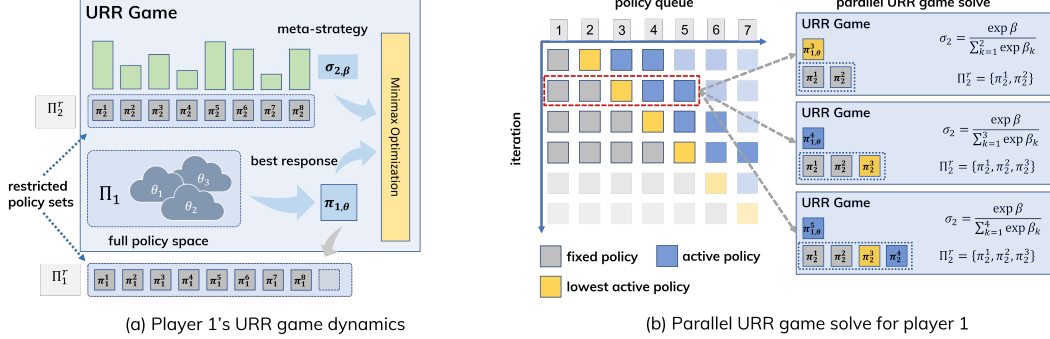(b) Parallel URR game solve for player 1

Figure 1: An overview of EPSRO. (a) EPSRO runs in a loop that learns best responses by solving URR games for each player, then expands the restricted policy sets with these $\pi_{i,\theta}$ at each epoch. (b) At each epoch, EPSRO runs multiple URR game solves for each active best response $\pi_{i,\theta}^j$ in parallel (Algorithm 4), where $j$ is the level. In each URR game, $\pi_{i,\theta}^j$ plays against $\Pi_{-i}^{r,k}$ where $k = 1, \ldots, j - 1$.

It is desirable for an efficient method to all of these problems. Our key insight is that the computation of meta-strategies can be free from simulations, and the learning of best responses should be toward monotonic expansion on restricted policy sets. As for the learning of best responses, there have several proposals for building an *unrestricted-restricted game* (URR game, Section 3.1) to learn a generalized best response as a Nash equilibrium to the opponent's restricted policy set [40, 16]. However, these existing works focus on tabular cases in which the policy space is limited.

In this paper, we introduce Efficient PSRO (EPSRO) based on the URR to save the simulation cost. EPSRO does not require simulations to compute meta-strategies beforehand. Instead, at each epoch, EPSRO solves URR games to learn a best response and an opponent meta-strategy, which is an approximate NE. In addition, we prove that the learned best responses are guaranteed to expand the restricted policy sets in a non-degenerate manner, improving the exploration efficiency. As for the URR solve, it is built on top of no-regret optimization [10]. Moreover, we propose an efficient warm-start technique for the regret optimization to save the re-training cost caused by the changed length of meta-strategies. We analyze the algorithm performance and give a regret bound of $\mathcal{O}(\sqrt{T \log [(k^2 + k)/2]})$, where $k$ is the size of the final restricted policy set. Most importantly, we introduce a pipeline URR solver to make the best response learning parallelizable, which further improves the exploration efficiency. The demonstration shows that EPSRO substantially improves the training efficiency and achieves better performance than existing PSRO-based methods in high-dimensional matrix games, poker games, and multi-agent gathering tasks.

## 2 Preliminaries

**Two-player Normal-form Games.** A two-player normal-form game [13] is a tuple $(\Pi, U^\Pi)$, where $\Pi = (\Pi_1, \Pi_2)$ and $U^\Pi = (U^{\Pi_1}, U^{\Pi_2})$ are the tuple of policy sets and the tuple of payoff tables, respectively. Formally, $\forall i \in \{1, 2\}, U^{\Pi_i} : \Pi \to \mathbb{R}^{|\Pi_1| \times |\Pi_2|}$, in which each item represents the utility of a joint policy. Players in the game try to maximize their own expected utility by sampling policy from a mixture (distribution) $\sigma_i$ over their policy sets, where $\forall i \in \{1, 2\}, \sigma_i \in \Delta(\Pi_i)$. For the sake of convenience, we use $-i$ to denote the other agent except for player $i$ in the following content. A best response to a mixed-strategy $\sigma_{-i}$ is defined as a strategy that has highest utility. It can be expressed as $\mathbf{BR}(\sigma_{-i}) = \arg\max_{\sigma_i'} u_i(\sigma_i', \sigma_{-i})$, where $u_i(\cdot, \cdot)$ represents the utility function of player $i$ for a given joint policy.

**Policy Space Response Oracles (PSRO)** Double Oracle (DO) methods [27, 20, 25] provide an iterative mechanism for finding a Nash equilibrium approximation in normal form games. These algorithms work in expanding a restricted policy set $\Pi_i^r$ for each player iteratively. At each epoch, a Nash equilibrium $\sigma = (\sigma_i, \sigma_{-i})$ is computed for a restricted game which is formed by a tuple of restricted policy sets $\Pi^r = (\Pi_i^r, \Pi_{-i}^r)$. Then, a best response to this Nash equilibrium for each player $i$ is computed and added to its restricted policy set $\Pi_i^r = \Pi_i^r \cup \{\mathbf{BR}(\sigma_{-i})\}$. PSRO is a generalization of DO where the restricted game's choice is a policy rather than an action. At each epoch, PSRO learns an approximate best response to an Nash equilibrium via the oracles (e.g., reinforcement learning algorithms). There are many different solvers for the computation of Nash equilibrium, such as $\alpha$-rank [28], PRD [19] or some linear programming methods [31]. In practice, PSRO seeks an approximation of Nash equilibrium, which is at a level of precision $\epsilon \geq 0$ [1]. To evaluate the equality of approximation, we use $\text{NASHCONV}(\sigma) = \sum_i u_i\left(\mathbf{BR}_i(\sigma_{-i}), \sigma_{-i}\right) - u_i(\sigma)$ to compute the exploitability of $\sigma$ to an oracle $\{\mathbf{BR}(\sigma_{-i})\}$ [18]. $\sigma$ is an exact Nash equilibrium if $\text{NASHCONV} = 0$.

We summarize the pseudo code of PSRO in Algorithm 1. At each epoch, PSRO requires simulations to compute the missing items in $U^{\Pi^r}$ after the learning of best responses, which causes an expensive computing cost. In general, the amount of simulations grows geometrically as $\mathcal{O}(M \cdot |\Pi_i^r|)$, where $|\Pi_i^r|$ and $M$ denote the size of restricted policy set and the number of simulations for each missing item, respectively. To learn approximate best responses, PSRO usually runs nested reinforcement learning algorithms. However, such a procedure is data-thirsty and has no guarantee to find a high-quality best response to bring higher payoffs for a restricted policy set, especially in the case of complex tasks.

---

**Algorithm 1:** VANILLA PSRO

**Input:** initial restricted policy sets $\Pi^r = (\Pi_1^r, \Pi_2^r)$
/* can be saved via URR games                */
**Input:** empty payoff table $U^{\Pi^r}$
**Input:** meta-strategies $\sigma_i \sim \text{UNIFORM}(\Pi_i^r)$
1 **while** *not terminated* **do**
2   **for** *player* $i \in \{1, 2\}$ **do**
3     **for** *many episodes* **do**
4       Train best response $\pi_{i,\theta}$ against $\pi_{-i} \sim \sigma_{-i}$
5     $\Pi_i^r = \Pi_i^r \cup \{\pi_i, \theta\}$
    /* can be saved via URR games                */
6   Run simulations to compute missing entries in $U^{\Pi^r}$
7   Compute a meta-strategy $\sigma$ from $U^{\Pi^r}$
**Output:** current meta-strategy $\sigma_i$ for player $i$

**Algorithm 2:** SIMPLIFIED PSRO WITH URR GAMES

**Input:** initial restricted policy sets $\Pi^r = (\Pi_1^r, \Pi_2^r)$
1 **while** *not terminated* **do**
2   **for** *player* $i \in \{1, 2\}$ **do**
3     Random initialize a best response $\pi_{i,\theta}$
4     $(\pi_{i,\theta}, \sigma_{-i,\beta}) = \text{SOLVEURR}(\pi_{i,\theta}, \Pi_{-i}^r)$
5   $\Pi_i^r = \Pi_i^r \cup \{\pi_{i,\theta}\}$ for $i \in \{1, 2\}$
**Output:** current meta-strategy $\sigma_i$ for player $i$

---

## 3 EPSRO: Efficient PSRO

For EPSRO, the keys to improve its efficiency include two aspects: (1) eliminating simulations for computing meta-strategies to improve the computing efficiency; (2) finding high-quality best responses to improve the exploration efficiency. We developed EPSRO on top of URR games (Section 3.1) for the first aspect. As for the exploration efficiency, it indicates the efficiency of expanding restricted policy sets. In general, higher exploration efficiency means the algorithm can express a restricted policy space with a smaller policy set than other methods. Thus, the quality of learned best responses is vital to the exploration. To tackle this problem, we propose an efficient algorithm to solve URR games in parallelism (Section 3.2 and 3.3). We summarize the pseudo code of Efficient PSRO (EPSRO) in Algorithm 4 and give its overview in Figure 1.

### 3.1 Modeling EPSRO as URR Games

Saving the computational cost of simulation is crucial to optimizing the efficiency of PSRO methods. In EPSRO, we achieve that by developing a simulation-free mechanism for meta-strategy learning and policy space expansion. Specifically, the meta-strategies are derived from direct interaction with best responses $\mathbf{BR}(\sigma_{-i}) \in \Delta(\Pi_i)$ instead of a simulation-based $U^{\Pi^r}$, and restricted policy sets are

expanded with these **BR**s. We further model the interaction as an *unrestricted-restricted game* below, which can be regarded as a parameterized extension of the tabular case in [40].

**Definition 3.1.** An *unrestricted-restricted* (URR) game for player $i$ is a tuple of full policy set $\Pi_i$ and restricted policy set $\Pi_{-i}^r$, i.e. $(\Pi_i, \Pi_{-i}^r)$. In this game, the player $i$ models its policy as a function parameterized by $\theta$, i.e. $\pi_{i,\theta} \in \Delta(\Pi_i)$. For each interaction, it plays against an opponent's policy $\pi_{-i}$ sampled from $\sigma_{-i,\beta} \in \Delta(\Pi_{-i}^r)$, where $\sigma_{-i,\beta}$ is a meta-strategy parameterized by $\beta$. $(\pi_{i,\theta}^\star, \sigma_{-i,\beta}^\star)$ is a Nash equilibrium if

$$\pi_{i,\theta}^\star = \mathbf{BR}(\sigma_{-i,\beta}^\star), \text{ and } \sigma_{-i,\beta}^\star = \mathbf{BR}(\pi_{i,\theta}^\star).$$

As described in Definition 3.1, the learning of best responses doesn't require fixed meta-strategies from a restricted game. Therefore, URR games save the computational cost of running simulations to construct the $U^{\Pi^r}$. Algorithm 2 lists the pseudo-code of URR-based PSRO for the comparison with the vanilla PSRO (Algorithm 1).

Though URR can save the simulation costs, we need to ensure its policy set expansion is non-degenerate since the final $\Pi^r$ should be an approximation of complete policy space $\Pi$. To evaluate how closely the restricted policy set is to $\Pi$, Balduzzi et al. [2] introduce the **gamescape**. Though the original cases are self-play, we can naturally bring this concept to URR games.

**Definition 3.2** (*URR Gamescape, derived from [2]*). Given an URR game $(\Pi_i, \Pi_{-i}^r)$ (row player $i$, column player $-i$) with payoff matrix $\mathbf{U}$, the corresponding empirical gamescape (EGS) is $\mathcal{G} := \{convex\ mixture\ of\ columns\ of\ \mathbf{U}\}$.

Conceptually, the non-degenerate or monotonic policy set expansion means that the learned best responses should not be in existing $\mathcal{G}$, so that the algorithm can expand the restricted policy set to approach the complete policy space. We investigate the policy set expansion of URR games in the following theorem.

**Theorem 3.3** (*Monotonic Policy Space Expanding*). *For any given epoch $e$ and $e+1$, let $(\pi_i^e, \sigma_{-i}^e)$ and $(\pi_i^{e+1}, \sigma_{-i}^{e+1})$ be Nash equilibrium of $\mathbf{URR}_i^e$ and $\mathbf{URR}_i^{e+1}$, respectively, where $\pi_i^e, \pi_i^{e+1} \in \Pi_i$, $\sigma_{-i}^e \in \Delta_{\Pi_{-i}^r}^e$ and $\sigma_{-i}^{e+1} \in \Delta_{\Pi_{-i}^r}^{e+1}$. The utilities of $\pi_i^e$ against opponent strategies $\sigma_{-i}^e$ and $\sigma_{-i}^{e+1}$ satisfies*

$$u_i(\pi_i^e, \sigma_{-i}^{\mathbf{e}}) - u_i(\pi_i^e, \sigma_{-i}^{\mathbf{e+1}}) \geq 0, \tag{1}$$

*where $\Delta_{\Pi_{-i}^r}^e$ indicates $\Delta(\Pi_{-i}^{r,e})$. Especially, $u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^e, \sigma_{-i}^{e+1}) > 0$ indicates there is a strictly policy space expanding at $e+1$, i.e., $\pi_{-i}^{e+1} \in \Pi_{-i}^{r,e+1} \setminus \Pi_{-i}^{r,e}$. (See Appendix B.1)*

Theorem 3.3 shows that EPSRO can achieve a monotonic performance improvement and policy space expansion with URR games. We further investigate that EPSRO has higher exploration efficiency than the naive PSRO (Section 4) and give a Proposition as follows.

**Proposition 3.4.** *EPSRO has higher exploration efficiency than PSRO. (See Appendix B.3)*

### 3.2 Solving URR Games

We've built our EPSRO with URR games and gave analysis for its policy set expansion, but how to solve the URR games is still a question. As solving Nash equilibrium in large-scale games is difficult, so we seek a learning-based method. We propose a multi-agent learning algorithm with the corresponding pseudo-code listed in Algorithm 3, which characterizes a procedure to train a best response and opponent meta-strategy for each player. Specifically, we learn the best response with a reinforcement learning algorithm while learning the meta-strategy with an online no-regret method.

---

**Algorithm 3:** SOLVEURR

**Input:** URR game $(\Pi_i, \Pi_{-i}^r)$, BR $\pi_{i,\theta} \sim \Delta(\Pi_i)$
**Input:** meta-strategies $\sigma_{-i,\beta} \sim$ UNIFORM$(\Pi_{-i}^r)$
1 **while** *not terminated* **do**
2      Train best response $\pi_{i,\theta'} \leftarrow \pi_{i,\theta}$ against $\pi_{-i} \sim \sigma_{-i,\beta}$ with reinforcement learning step
3      Update $\sigma_{-i,\beta'} \leftarrow \sigma_{-i,\beta}$ against $\pi_{i,\theta}$ by following Algorithm 5 in Appendix A
4      $\pi_{i,\theta} \leftarrow \pi_{i,\theta'}, \sigma_{-i,\beta} \leftarrow \sigma_{-i,\beta'}$
**Output:** an approximate Nash $(\pi_{i,\theta}^\star, \sigma_{-i,\beta}^\star)$

---

In Algorithm 3, the opponent meta-strategy $\sigma_{-i}$ is represented as a Boltzman distribution, which is parameterized by a vector $\beta = \left[\beta_1, \beta_2, \ldots, \beta_{|\Pi^r_{-i}|}\right]$. Thus, each support of $\sigma_{-i}$ could be expressed as $\sigma_{-i}(j) = \exp \beta_j / \sum_{i=1}^n \exp \beta_i$. We update $\sigma_{-i}$ by following Algorithm 5 in Appendix C. As for the best response $\pi_i$, it is a neural network trained with off-policy reinforcement learning [33].

A feasible method to quantify the learning performance of Algorithm 3 is to calculate the regret to the oracle payoff. In this paper, we use no-regret algorithms [4, 10] to analyze EPSRO's algorithm performance in two-player zero-sum games. Under this framework, a learning algorithm could approximate the NE asymptotically by playing the same game repeatedly. A well-known no-regret algorithm among them is Multiplicative Weights Update (MWU) [12], which updates strategy by considering the averaging loss along the learning horizon and then achieves no-regret as the learning horizon towards infinite. Though the update of EPSRO's meta-strategies does not exploit the loss directly like MWU, it follows MWU by exploiting the computed gradients. We further explain it in Lemma B.5.

**Definition 3.5.** Considering a sequence of mixed strategies for player $i$ as $\pi_1, \pi_2, \ldots$, an algorithm of $-i$ that generates a sequence of mixed strategies $\sigma_1, \sigma_2, \ldots$ is called a no-regret algorithm if we have: $\lim_{T \to \infty} R_T / T = 0$, $R_T = \max_{\sigma \in \Delta_{\Pi_{-i}}} \sum_{t=1}^T (\pi_t^\top U \sigma - \pi_t^T U \sigma_t)$.

**Warm-Start Learning.** In EPSRO, a critical problem for the meta-strategy optimization is that the length of meta-strategy changes as the policy set expands. Therefore, we need to learn the meta-strategy with a new $\beta$ at each epoch. However, if we start the learning from scratch, the cost of re-training grows more and more expensive as the policy set expands, resulting in a slow convergence rate. To tackle this issue, we propose a warm-start technique that enables the meta-strategy starts from a non-trivial initialization. Since the warm-start aims to save time to achieve a minimal regret when $\Pi^r$ changes, the key is to correctly initialize the regrets instead of only starting the strategy. Moreover, a wrong initialization of the regrets will result in huge regrets in subsequent iterations, which is no better than starting from scratch. As pointed in [6], a feasible warm-start should be a substitute strategy that does not violate the regret bound and the approximate NE of the last epoch.

**Theorem 3.6** (*Theorem 1 in [6]*). *In a two-player zero-sum game, if $\frac{R_i^T}{T} \leq \epsilon_i$ for both player $i \in \{1, 2\}$, then $(\bar{\pi}_i, \bar{\sigma}_{-i})$ is a $(\epsilon_1 + \epsilon_2)$-equilibrium, where $\bar{\sigma}_{-i} = \sum_{t=1}^T \langle \sigma_{-i,t}, l_{-i,t} \rangle / \sum_{t=1}^T T l_{-i,t}$, $\bar{\pi}_i = \sum_{t=1}^T \langle \pi_{i,t}, l_{i,t} \rangle / \sum_{t=1}^T T l_{i,t}$, $R_i^T$ the summation of regrets of $T$ iterations, $l_t$ the loss vector.*

Theorem 3.6 shows that if we use a regret-based average $\bar{\sigma}_{-i}$ (or $\bar{\pi}_i$) as the substitute of a sequence of $\sigma_{-i,t}$ (or $\pi_{i,t}$), the substitute can still hold the equilibrium. In case of EPSRO, as the restricted policy set grows from $\Pi^{r,e}_{-i}$ to $\Pi^{r,e+1}_{-i}$, we need to investigate whether there is a feasible substitute $\bar{\sigma}'_{-i} \in \Delta^{e+1}_{\Pi^r_{-i}}$ to $\bar{\sigma}_{-i} \in \Delta^e_{\Pi^r_{-i}}$ follows the regret guarantee of epoch $e$.

**Theorem 3.7** (*See Appendix B.6*). *Suppose a substitute policy of $\bar{\sigma}_{-i} \in \Delta^e_{\Pi^r_{-i}}$ is $\bar{\sigma}'_{-i} \in \Delta^{e+1}_{\Pi^r_{-i}}$, and it satisfies $u_i^{e+1}(\bar{\pi}_i, \bar{\sigma}'_{-i}) = u_i^e(\bar{\pi}_i, \bar{\sigma}_{-i})$, we have*

$$\max_{\sigma_{-i} \in \Delta^{e+1}_{\Pi^r_{-i}}} \sum_{t=1}^T \left( u_{-i}^{e+1}(\pi_i^t, \sigma_{-i}) - u_{-i}^{e+1}(\pi_i^t, \bar{\sigma}'_{-i}) \right) \leq \epsilon_{-i},$$

*where $\epsilon_{-i}$ is the regret bound of epoch $e + 1$, $\pi_i^t \in \Delta(\Pi_i)$.*

Theorem 3.7 shows that there exists a substitute $\bar{\sigma}'_{-i} \in \Delta^{e+1}_{\Pi^r_{-i}}$ holds Theorem 3.6 also holds the regret bound of epoch $e + 1$. Thus, once we compute such a $\bar{\sigma}'_{-i}$ that satisfies $u_i(\bar{\pi}_i, \bar{\sigma}'_{-i}) = u_i(\bar{\pi}_i, \bar{\sigma}_{-i})$, then it is a feasible warm-start at the next epoch to save the training time.

**Lemma 3.8** (*See Appendix B.7*). *Let $k = |\Pi^{r,e+1}_{-i}|$, $\bar{\sigma}'_{-i}$ be parameterized by $\beta_{-i} = [\beta_{-i,1}, \beta_{-i,2}, \ldots, \beta_{-i,k}]$, $\bar{\sigma}'_{-i}(k)$ the $k$-th item of $\bar{\sigma}'_{-i}$, $x = [\bar{\sigma}'_{-i}(1), \ldots, \bar{\sigma}'_{-i}(k-1)]^T$, $\bar{l}^e_{-i}$ is $-i$'s average loss vector to $\bar{\pi}_i$ at epoch $e$. Then a feasible initial of $\beta^{e+1}_{-i}$ could satisfy*

$$\beta^{e+1}_{-i} = \arg\min_{\beta_{-i}} \| (x - \bar{\sigma}_{-i})^\top \bar{l}^e_{-i} - \bar{\sigma}'_{-i}(k) u_{-i}(\bar{\pi}_i, \pi^{e+1}_{-i}) \|_2 . \tag{2}$$

5

Lemma 3.8 computes a $\beta_{-i}^{e+1}$ for the initialization of $\sigma_{-i}^{e+1}$ to satisfy Theorem 3.6 and 3.7. In addition, the best response $\pi_{i,\theta}$ also follows these conditions by continuous training instead of restarting parameters at each epoch. We now present the regret bound of EPSRO as follows.

**Theorem 3.9** (*Regret Bound of EPSRO*). *Let $l_1, l_2, \ldots, l_T$ be a sequence of loss vectors player by an opponent, and $\langle \cdot, \cdot \rangle$ be the dot product, then EPSRO is a no-regret algorithm with (See Appendix B.8)*

$$\frac{1}{T} \left( \sum_{t=1}^{T} \langle \sigma_t, l_t \rangle - \min_{\sigma \in \Delta(\Pi_t^r)} \sum_{t=1}^{T} \langle \sigma, l_t \rangle \right) \leq \frac{\sqrt{\log[(k+1)k/2]}}{\sqrt{2T}}, \text{ where } k \text{ is the size of } \Pi^r.$$

Theorem 3.9 shows that EPSRO is no-regret when the policy set is finite. Though some complex games have continuous policy space, the number of effective policies is finite. We further give the convergence rate of EPSRO as follows.

**Theorem 3.10** (*Convergence Rate of EPSRO*). *Let $k$, $N$ denote the size of restricted policy sets $\Pi_{-i}^r$ and $\Pi_i$. Then the learning of Algorithm 3 will converge to the Nash equilibrium with the rate (See Appendix B.9):*

$$\epsilon_T = \sqrt{\frac{\log[(k+1)k/2]}{2T}} + \sqrt{\frac{\log[(N+1)N/2]}{2T}}.$$

### 3.3   Pipeline URR Solver

While Algorithm 3 is clear, the reinforcement learning step can take a long time to converge to a good response, especially in complex games. To overcome this training problem, we introduce a parallel solution to further improve EPSRO's training efficiency. Inspired by Pipeline-PSRO (P-PSRO) [24], we schedule multiple asynchronous training procedures to train a best response policy for each player. Figure 1 illustrates an example dynamics, it able to scale up the *minimax optimization* with convergence guarantee by maintaining a hierarchical pipeline of reinforcement learning

---

**Algorithm 4:** EFFICIENT PSRO (EPSRO)

**Input:** inital restricted policy sets $\Pi^r = (\Pi_1^r, \Pi_2^r)$

**1  while** *not terminated* **do**
**2**      **for** *all player $i \in \{1, 2\}$ in parallel* **do**
**3**          **for** *loop all active best response $\pi_i^j \in \Pi_i^r$* **do**
**4**              **for** *all $\Pi_{-i}^{r,<j}$ in parallel* **do**
**5**                  $\pi_{i,\theta}^j, \sigma_{-i,\beta}^{<j} = \text{SOLVEURR}(\pi_{i,\theta}^j, \Pi_{-i}^{r,<j})$
**6**                  **if** *the lowest $\pi_{i,\theta}^j$ meets stops cond.* **then**
**7**                      Set it to fixed and $\Pi_i^r := \Pi_i^r \cup \{\pi_i^j\}$
**8**                      Generate a new active policy

**Output:** current meta-strategies $\sigma = (\sigma_1, \sigma_2)$

---

policies as P-PSRO (Proposition 3.2 in [24]). We finally give EPSRO's pseudo-code with parallel URR solve in Algorithm 4. Each player $i$ in parallel maintains a queue of ordered policies with two classes of training policies: fixed policies at low levels and active policies at high levels. Each active policy $\pi_i^j$ at level $j$ learns against the opponent restricted policy set $\Pi_{-i}^{r,<j} \in \{\Pi_{-i}^{r,k} | k = 1, \ldots, j-1\}$ which is composed of active and fixed policies lower than $j$. Once a lowest active policy meets the stop condition (e.g. number of training episodes), it will be fixed to expand the corresponding restricted policy set, and a new active policy will be add into the queue with the highest level. We argue that parallel training will increase exploration efficiency since each best response plays against multiple policy sets.

## 4   Experiments

We compare EPSRO with five algorithms, including Self-play [17], PSRO [19], Rectified PSRO (PSRO-rN) [2], Mixed-Oracles [32], and Pipeline PSRO (P-PSRO) [24]. The environments for the test are three classes of increasing difficulty games, i.e., matrix games, Poker games, and Multi-agent Gathering. We investigate the *exploration efficiency* and *computation efficiency* of algorithms in these experiments, also the performance on convergence. We use NASHCONV to evaluate the convergence quality in matrix games and Poker games and *cardinality of payoff matrix* [30] to evaluate the exploration efficiency. Since traversing the game tree of the Multi-agent Gathering game is too expensive, we set the policy sets $\Pi^{\text{PSRO}}$ produced by PSRO as the baseline to calculate the score of algorithms, which reflects the performance to some extent. The matrix games are designed for

the comparison of *exploration efficiency*. Especially the non-transitive mixture game (Section 4.1), which vividly characterizes the exploration dynamics of algorithms. Furthermore, to demonstrate the importance of parallel training for exploration efficiency, we remove the pipeline URR solver of EPSRO as NEPSRO in the matrix games. As for the *computation efficiency*, we investigate it from the number of samples and time consumption in Poker games and Multi-agent Gathering. Extra results and the pseudo-code of score calculation in Appendix F. All experiments were performed on a single machine with 64 CPUs, 256 G RAM, and 2 GeForce RTX 3090 GPUs.



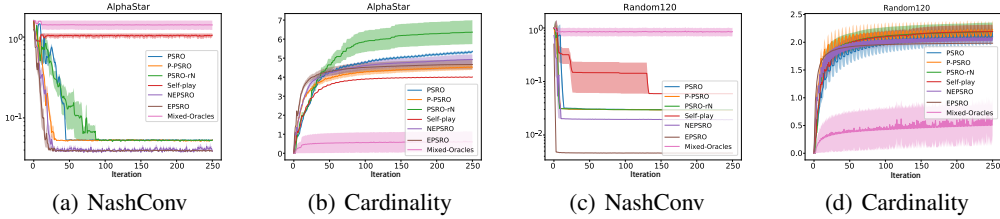(a) NashConv        (b) Cardinality        (c) NashConv        (d) Cardinality

Figure 2: Comparison of NASHCONV and cardinality. (a) and (b) show the NASHCONV and cardinality on AlphaStar matrix game, respectively; (c) and (d) show the NASHCONV and cardinality on a high-dimensional symmetric game ($n = 120$), respectively. Though EPSRO has lower cardinality than some other algorithms, it outperforms all baselines on the NASHCONV. We argue that a lower cardinality but better exploitability indicates that the algorithm has higher exploration efficiency since it achieves better performance with fewer policies. More results in Appendix F.2.
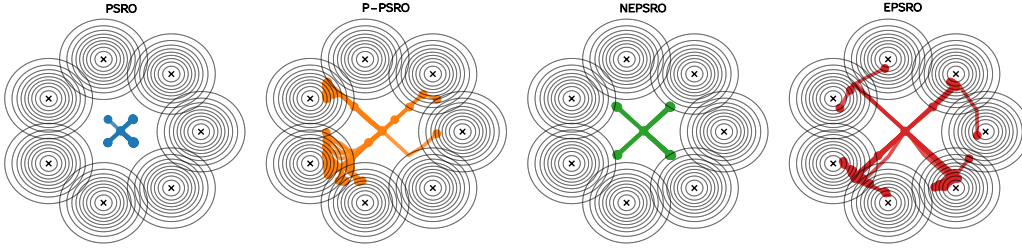


Figure 3: Exploration trajectories on *Non-transitive Mixture Games*. The more trajectories close to the centers of Gaussian, the higher the exploration efficiency of the algorithm. Our algorithms (EPSRO and NEPSRO) outperform all selected baselines. Especially the EPSRO, it explored all centers. More results in Appendix F.1.
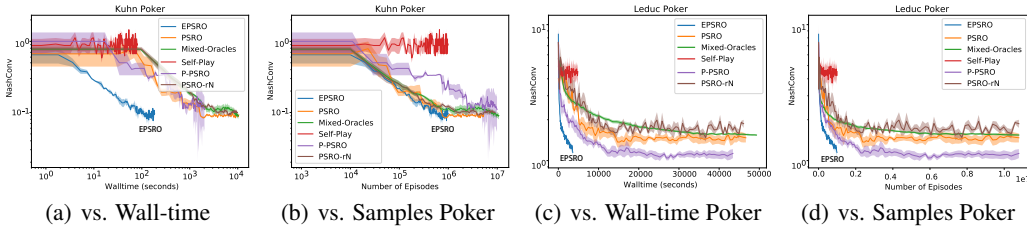


(a) vs. Wall-time        (b) vs. Samples Poker        (c) vs. Wall-time Poker        (d) vs. Samples Poker

Figure 4: NashConv on poker games. The number of samples for training each best response is set to $1e4$ episodes, and the number of simulations (if needed) for each joint policy is set to $1e3$ episodes. EPSRO performs a similar performance on the NASHCONV. For the computation efficiency, EPSRO achieves more than 50x seepup on wall-time; more than 10x sample efficiency than other algorithms.

## 4.1 Comparison of Exploration Efficiency

The non-transitive game for the comparison of *exploration efficiency* is a zero-sum two-player game consisting of 7 equally-distanced Gaussian humps on the 2D plane. In this game, each player chooses a point in the 2D plane as its decision, which is transformed into a 7-dimensional vector $\pi_i$ with each

coordinate being the density in the corresponding Gaussian distribution. The payoff of the game is given by $\phi_i(\pi_i, \pi_{-i}) = \pi_i^T \mathbf{S} \pi_{-i} + \mathbf{1}^T(\pi_i - \pi_{-i})$. In this game, the optimal strategy should stay close to the center of the Gaussian and explore all the Gaussian distributions equally. We train best response policies in 50 epochs for each algorithm. As presented in Figure 3, EPSRO successfully explores all the centers and shows higher *exploration efficiency* than other baselines. Though the NEPSRO fails to achieve to any centers, its explored policy space is larger than most algorithms.

## 4.2 High-dimensional Matrix Games

The high-dimensional matrix games introduced here include two classes. One is a symmetric matrix game generated in a high-dimensional uniform distribution, and another is an empirical payoff matrix corresponding to 888 reinforcement learning policies in AlphaStar [35]. We demonstrate the comparison of these games on the performance when dealing with games that have high-dimensional policy spaces.

**Random Symmetric Matrix Game.** [24] introduce the games to investigate the performance of PSRO-based methods in high-dimensional symmetric games (SymGame). In this experiment, we generated random symmetric zero-sum matrices with different dimension $n$. For a given matrix, elements in the upper triangle are distributed uniformly: $\forall i < j \leq n, a_{i,j} \sim \text{UNIFORM}(-1, 1)$ and for the lower triangle, the elements are set to be the negative of its diagonal counterpart: $\forall j < i \leq n, a_{i,j} = -a_{j,i}$. The diagonal elements are equal to zero: $a_i, i = 0$. The matrix defines the utility of two pure strategies to the row player. In these experiments, we train a strategy $\pi$ as a best response that plays against another strategy $\hat{\pi}$, it is updated by a learning rate $r$ multiplied by the best response to that strategy: $\pi' = r\mathbf{BR}(\hat{\pi}) + (1 - r)\pi$. Figure 2 shows the results for $n = 120$. We report both NASHCONV and cardinality. The results show that EPSRO and NEPSRO achieve a faster convergence rate and the lowest NASHCONV than all of the other algorithms. Though they do not achieve the highest cardinality, we argue that there is a tradeoff between convergence and exploration, and EPSRO performs a better balance between them. It is worth mentioning that the Mixed-Oracle fails to seek a meta-strategy that has a smaller distance to the Nash equilibrium, even worse than Self-Play. We argue that the policy distills of Mixed-Oracles may decrease the exploration efficiency, especially in such a high-dimensional policy space.

**AlphaStar Empirical Game.** The AlphaStar Matrix Game is derived from solving a complex real-world game StraCraftII [8], which involves 888 reinforcement learning policies. We test the *exploration efficiency* and the convergence quality of our method for solving such empirical games. Similar to the results in the random symmetric matrix game, our algorithm performs a faster convergence rate and lower NASHCONV than other algorithms, while the Mixed-Oracle still fails to explore new policies to expand its policy sets (Figure 2).

## 4.3 Poker Games

Poker is a common benchmark in multi-agent decision tasks. In this paper, we introduce two simplified forms of poker games for experiments, i.e., Kuhn Poker and Leduc Poker [41, 5]. These poker tasks model zero-sum two-player imperfect-information games, in which each player shows uncertainty about the game rules and the state of other players. Similar to Poker, where each player in these games chooses to raise/call/fold through rounds of betting. We investigate the sample efficiency and performance of EPSRO in this experiment. The training for each



Figure 5: The average score of final policy set that plays against $\Pi^{\text{PSRO}}$.

algorithm is set to learn 100 policies. The number of samples for training each policy is set to 10000 episodes, and the number of simulations for each joint policy is set to 1000 episodes. So the total number of samples for training a PSRO algorithm with simulations achieves $1.1 \times 1e7$. Figure 4 presents the results of NACHCONV w.r.t to wall-time and the number of samples. Since no simulations and parallel best response learning, EPSRO substantially achieves high training efficiency. Specifically, for the wall-time, EPSRO has more than 50x speedup than other PSRO methods; for the sample efficiency, EPSRO has more than 10x than other non-parallelized PSRO methods.
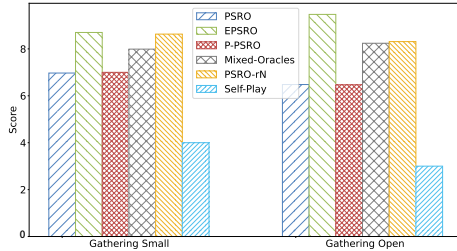
### 4.4 Multi-agent Gathering

We further investigate the capability of EPSRO to handle more complex multi-agent tasks in multi-agent gathering environments (MAG)[1]. MAG is an endless environment whose horizon length could be infinite. In our experiments, we limit the horizon of each episode to 100. In MAG, the goal of each agent is to collect as many as "apples". The apples regrow at a rate dependent upon the configuration of the uncollected nearby apples. In this case, the more nearby apples, the higher the regrowth rate of the apples. Naturally, this presents a dilemma for the players: each wants to pick as many apples as possible. However, if they over-harvest the throughput of apples diminishes, potentially falling to zero. Figure 5 shows the confrontation results of each algorithm's final policy set $\Pi^{\text{TEST}}$ to an evaluation policy set $\Pi^{\text{PSRO}}$. We calculate the score as $\sigma^{\text{TEST}} M^{\text{TEST}} \left[\sigma^{\text{PSRO}}\right]^{T}$, where $M^{\text{TEST}}$ is the empirical payoff matrix and $\sigma$ is a learned meta-strategy. The higher the score, the better the performance of the corresponding algorithm. Additionally, we report the curve of the score of intermediate policy sets during the learning process in Appendix F.3.

## 5  Related Work

In large normal-form games, it is difficult to directly compute the approximation of Nash equilibrium. Policy Space Response Oracles (PSRO) [19] provides an iterative solution to solve this problem. There are many variants of PSRO focus to improve the convergence rate or training efficiency. A straightforward idea is to utilize parallelism to improve the training efficiency. For instance, DCH [19] parallelizes PSRO by training multiple RL policies, each against the meta Nash equilibrium below it in the hierarchy. A problem with this method is that the number of policies should be set beforehand. However, it is difficult to figure out how many policies does it require to solve a game in practice. [24] proposed a similar solution, i.e. P-PSRO, to solve this problem, which inherits the hierarchical parallelism training but has no need to preset the number of training policies. Another parallelism variant is Rectified PSRO [2], but it has been proved not converge in all symmetric zero-sum games [24].

Another factor to the efficiency is the exploration efficiency. Specifically, more diverse the learned best responses, closer the restricted sets to original policy sets [29, 37]. Thus, PSRO-based methods converge in smaller size of restricted policy sets and less time consumption. However, it is difficult to discover an exact best response, because the underlying policy training has no guarantee to find an ideal policy as the best response without extra conditions. As a solution, improving the diversity of restricted policy set has been regarded as a reasonable way to solve this problem. Among the existing work, DPP [29] utilizes the expected cardinality to measure the diversity of policy set. [23] proposed a method that unify both behavior and reward distance to measure the diversity. Since the learning of best response need to play against a mixture of opponent policies. Many of existing work demonstrate by sampling opponent policies from this mixture for each episode. Compared to playing against a single policy, [32] claimed that such a mechanism brings stochasticity on the opponent and forgets previous experiences, making algorithms slow to converge. Thus, the authors proposed a method that distills [9] the opponent mixture as a single policy via Q-mixing [11].

Concurrently to our work, [26] proposed a similar work named AODO. Their work differs from ours in the following ways: (1) We focus on improving the training efficiency of PSRO-based methods, providing monotone improvement and convergence analysis for EPSRO, while AODO focuses on making a guarantee to monotonic decrease the NASHCONV [18]; (2) EPSRO introduces parallelized best response training while AODO executes in singleton; (3) We proposed a warm-start technique to tackle the re-training problem while AODO starts from scratch at each epoch; (4) We demonstrate the experiments with 5 baselines on both high-dimensional matrix games, poker games, and non-trivial multi-agent gathering tasks while AODO considers only one baseline and fewer environments.

## 6  Conclusions

In this work, we propose a parallel algorithm EPSRO to improve the training efficiency of PSRO. The demonstration results show that EPSRO achieves higher computation efficiency and exploration efficiency than existing works. The improvements of EPSRO benefit from learning best responses

---

[1] https://github.com/HumanCompatibleAI/multi-agent

against the whole opponent restricted policy set and cooperating with parallelized training. In addition, a warm-start technique to reduce the re-training cost also makes our algorithm perform a higher convergence rate. However, EPSRO is limited to handling the two-player cases because there will be a divergence in selecting meta-strategies for more players involved. In future work, we would like to seek a method to solve this problem and generalize EPSRO to multi-player cases.

# References

[1] Haris Aziz. Multiagent systems: algorithmic, game-theoretic, and logical foundations by y. shoham and k. leyton-brown cambridge university press, 2008. *ACM Sigact News*, 41(1):34–37, 2010.

[2] David Balduzzi, Marta Garnelo, Yoram Bachrach, Wojciech Czarnecki, Julien Perolat, Max Jaderberg, and Thore Graepel. Open-ended learning in symmetric zero-sum games. In *International Conference on Machine Learning*, pages 434–443. PMLR, 2019.

[3] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. 2019.

[4] Michael Bowling. Convergence and no-regret in multiagent learning. *Advances in neural information processing systems*, 17, 2004.

[5] Michael Bowling, Neil Burch, Michael Johanson, and Oskari Tammelin. Heads-up limit hold'em poker is solved. *Science*, 347(6218):145–149, 2015.

[6] Noam Brown and Tuomas Sandholm. Strategy-based warm starting for regret minimization in games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.

[7] Shicong Cen, Fan Chen, and Yuejie Chi. Independent natural policy gradient methods for potential games: Finite-time global convergence with entropy regularization. *arXiv preprint arXiv:2204.05466*, 2022.

[8] Wojciech M Czarnecki, Gauthier Gidel, Brendan Tracey, Karl Tuyls, Shayegan Omidshafiei, David Balduzzi, and Max Jaderberg. Real world games look like spinning tops. *Advances in Neural Information Processing Systems*, 33, 2020.

[9] Wojciech M Czarnecki, Razvan Pascanu, Simon Osindero, Siddhant Jayakumar, Grzegorz Swirszcz, and Max Jaderberg. Distilling policy distillation. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1331–1340. PMLR, 2019.

[10] Constantinos Daskalakis, Alan Deckelbaum, and Anthony Kim. Near-optimal no-regret algorithms for zero-sum games. In *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pages 235–254. SIAM, 2011.

[11] Vasilii Davydov, Alexey Skrynnik, Konstantin Yakovlev, and Aleksandr Panov. Q-mixing network for multi-agent pathfinding in partially observable grid environments. In *Russian Conference on Artificial Intelligence*, pages 169–179. Springer, 2021.

[12] Yoav Freund and Robert E Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29(1-2):79–103, 1999.

[13] Drew Fudenberg and Jean Tirole. *Game theory*. MIT press, 1991.

[14] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, pages 1352–1361. PMLR, 2017.

[15] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018.

[16] Thomas Dueholm Hansen, Peter Bro Miltersen, and Troels Bjerre Sørensen. On range of skill. In *AAAI*, pages 277–282, 2008.

[17] Daniel Hernandez, Kevin Denamganaï, Yuan Gao, Peter York, Sam Devlin, Spyridon Samoth-rakis, and James Alfred Walker. A generalized framework for self-play training. In *2019 IEEE Conference on Games (CoG)*, pages 1–8. IEEE, 2019.

[18] Michael Johanson, Kevin Waugh, Michael Bowling, and Martin Zinkevich. Accelerating best response calculation in large extensive games. In *Twenty-second international joint conference on artificial intelligence*, 2011.

[19] Marc Lanctot, Vinicius Zambaldi, Audrūnas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent rein-forcement learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 4193–4206, 2017.

[20] Yaodong Yang Le Cong Dinh, Zheng Tian, Nicolas Perez Nieves, Oliver Slumbers, David Henry Mguni, Haitham Bou Ammar, and Jun Wang. Online double oracle. 2021.

[21] David S Leslie and Edmund J Collins. Generalised weakened fictitious play. *Games and Economic Behavior*, 56(2):285–298, 2006.

[22] Siqi Liu, Luke Marris, Daniel Hennes, Josh Merel, Nicolas Heess, and Thore Graepel. Neupl: Neural population learning. In *International Conference on Learning Representations*, 2021.

[23] Xiangyu Liu, Hangtian Jia, Ying Wen, Yaodong Yang, Yujing Hu, Yingfeng Chen, Changjie Fan, and Zhipeng Hu. Towards unifying behavioral and response diversity for open-ended learning in zero-sum games. *Advances in Neural Information Processing Systems*, 34, 2021.

[24] Stephen McAleer, JB Lanier, Roy Fox, and Pierre Baldi. Pipeline psro: A scalable approach for finding approximate nash equilibria in large games. *Advances in Neural Information Processing Systems*, 33:20238–20248, 2020.

[25] Stephen McAleer, John Banister Lanier, Kevin A Wang, Pierre Baldi, and Roy Fox. Xdo: A double oracle algorithm for extensive-form games. *Advances in Neural Information Processing Systems*, 34:23128–23139, 2021.

[26] Stephen McAleer, Kevin Wang, Marc Lanctot, John Lanier, Pierre Baldi, and Roy Fox. Anytime optimal psro for two-player zero-sum games, 2022.

[27] H Brendan McMahan, Geoffrey J Gordon, and Avrim Blum. Planning in the presence of cost functions controlled by an adversary. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 536–543, 2003.

[28] Shayegan Omidshafiei, Christos Papadimitriou, Georgios Piliouras, Karl Tuyls, Mark Rowland, Jean-Baptiste Lespiau, Wojciech M Czarnecki, Marc Lanctot, Julien Perolat, and Remi Munos. $\alpha$-rank: Multi-agent evaluation by evolution. *Scientific reports*, 9(1):1–29, 2019.

[29] Nicolas Perez-Nieves, Yaodong Yang, Oliver Slumbers, David H Mguni, Ying Wen, and Jun Wang. Modelling behavioural diversity for learning in open-ended games. In *International Conference on Machine Learning*, pages 8514–8524. PMLR, 2021.

[30] Jean-Charles Régin and Carla P Gomes. The cardinality matrix constraint. In *International Conference on Principles and Practice of Constraint Programming*, pages 572–587. Springer, 2004.

[31] Tuomas Sandholm, Andrew Gilpin, and Vincent Conitzer. Mixed-integer programming methods for finding nash equilibria. In *AAAI*, pages 495–501, 2005.

[32] Max Smith, Thomas Anthony, and Michael Wellman. Iterative empirical game solving via single policy best response. In *International Conference on Learning Representations*, 2020.

[33] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[34] Jan van den Brand. A deterministic linear program solver in current matrix multiplication time. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 259–278. SIAM, 2020.

[35] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Jun-young Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.

[36] Yongzhao Wang, Qiurui Ma, and Michael P Wellman. Evaluating strategy exploration in empirical game-theoretic analysis. *arXiv preprint arXiv:2105.10423*, 2021.

[37] Yaodong Yang, Jun Luo, Ying Wen, Oliver Slumbers, Daniel Graves, Haitham Bou Ammar, Jun Wang, and Matthew E Taylor. Diverse auto-curriculum is critical for successful real-world multiagent learning systems. In *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems*, pages 51–56, 2021.

[38] Yaodong Yang, Rui Luo, Minne Li, Ming Zhou, Weinan Zhang, and Jun Wang. Mean field multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 5571–5580. PMLR, 2018.

[39] Yaodong Yang, Rasul Tutunov, Phu Sakulwongtana, and Haitham Bou Ammar. $\alpha\alpha$-rank: Practically scaling $\alpha$-rank through stochastic optimisation. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, pages 1575–1583, 2020.

[40] Martin Zinkevich, Michael Bowling, and Neil Burch. A new algorithm for generating equilibria in massive zero-sum games. In *PROCEEDINGS OF THE NATIONAL CONFERENCE ON ARTIFICIAL INTELLIGENCE*, volume 22, page 788. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2007.

[41] Martin Zinkevich, Michael Johanson, Michael Bowling, and Carmelo Piccione. Regret minimization in games with incomplete information. *Advances in neural information processing systems*, 20:1729–1736, 2007.

# Appendices

## A  Algorithms

Algorithm 5 shows the pseudo-code of meta-strategy optimization.

---

**Algorithm 5:** DETERMINISTIC META STRATEGY OPTIMIZATION

**Input:** initialize the $\sigma_{-i,\beta}$ with Algorithm 6
**Input:** initialize an episode reword buffer $\mathcal{B}$; window size $L$; $K$ the length of meta-strategy $\sigma_{-i,\beta}$
**Input:** loss vector buffer $\mathcal{L}_i, \mathcal{L}_{-i} \in \mathbb{R}^K$; counters: $N_1 = 0, \dots, N_K = 0$

1  **while** *not terminated* **do**
2     **for** *many episodes* **do**
3        Sample $\pi^k_{-i} \sim \sigma_{-i,\beta}$ to play against $\pi_{i,\theta}$
4        Collect episode return $r^k$ into $\mathcal{B}$ and $N_k := N_k + 1$
5        **if** *the size of $\mathcal{B}$ equals to $L$* **then**
6           Calculate average observed returns: $\forall k = 1, \dots, K, \bar{r}^k = \frac{\sum_{l=1}^L \mathbf{1}_{j=k} r_L^j}{N_k}$
7           Compute gradients for $\beta_k = \arg\max_{\beta_k} \bar{r}^k$ and update $\sigma_{-i,\beta}$ with Lemma B.7
8           Collect loss vectors $l^k_{-i} = -\bar{r}^k, l^k_i = \bar{r}^k$ into $\mathcal{L}_i, \mathcal{L}_{-i}$, respectively
9           Reset buffer $\mathcal{B}$ and counters

**Output:** meta-strategy $\sigma_{-i}$ for player $-i$, loss vector buffers $\mathcal{L}_i, \mathcal{L}_{-i}$

---

We introduce the warm-start for opponent meta-strategy as follows. For each player $i$, the algorithm randomly initializes a vector $\beta$ with the length equals to the size of $\Pi^r_{-i}$. At each epoch $e + 1$, we slightly run a simulation procedure of $(\pi_i, \pi_{-i})$ to estimate the utility $u_{-i}(\pi_i, \pi_{-i})$ before applying the Lemma 3.8 to optimize $\beta$. Then, by cooperating with the average loss vector from last epoch $e$, we could start the optimization to compute a feasible $\beta^*$ that makes $\sigma_{-i}$ satisfy Theorem 3.6 and 3.7.

---

**Algorithm 6:** META-STRATEGY WARM-START

**Input:** newly learned policy support $\pi_{-i}$; best response $\pi_i$ from last epoch; error threshold $\tau \geq 0$
**Input:** initialize $\sigma_{-i}$ with $\beta \in \mathbb{R}^{|\Pi^r_{-i}|}$; loss vectors from last epoch $[l^1_{-i}, \dots, l^T_{-i}]$

1  Compute average loss vector $\bar{l}_{-i} = \left( l^1_{-i} + \cdots + l^T_{-i} \right) / T$
2  Estimate $u_{-i}(\pi_i, \pi_{-i})$ by running simulation for $(\pi_i, \pi_{-i})$
3  Compute the square error of utility as
   $\xi_2 = \| (x - \sigma_{-i})^\top \bar{l}_{-i} - \sigma_{-i}(k) u_{-i}(\pi_i, \pi_{-i}) \|_2 - \lambda H(\sigma_{-i})$
4  Update $\beta$ as $\beta := \beta - \nabla_\beta \xi_2$
5  **if** $\xi_2 \leq \tau$ **then**
6     Stop optimization and continue
7  **else**
8     Go back to step 3

**Output:** meta-strategy $\sigma_{-i}$

---

Algorithm 6 relies on the computation of $\xi_2$, we introduce its definition in Lemma B.7.

## B  Proofs

### A Proof of Theorem 3.3

**Theorem B.1** (*Monotonic Policy Space Expanding*)**.** *For any given epoch $e$ and $e + 1$, let $(\pi^e_i, \sigma^e_{-i})$ and $(\pi^{e+1}_i, \sigma^{e+1}_{-i})$ be Nash equilibrium of $\mathbf{URR}^e_i$ and $\mathbf{URR}^{e+1}_i$, respectively, where $\pi^e_i, \pi^{e+1}_i \in \Pi_i$, $\sigma^e_{-i} \in \Delta^e_{\Pi_{-i}}$ and $\sigma^{e+1}_{-i} \in \Delta^{e+1}_{\Pi^r_{-i}}$. The utilities of $\pi^e_i$ against opponent strategies $\sigma^e_{-i}$ and $\sigma^{e+1}_{-i}$*

*satisfies*

$$u_i(\pi_i^e, \sigma_{-\mathbf{i}}^{\mathbf{e}}) - u_i(\pi_i^e, \sigma_{-\mathbf{i}}^{\mathbf{e+1}}) \geq 0, \tag{3}$$

*where $\Delta_{\Pi_{-i}^r}^e$ indicates $\Delta(\Pi_{-i}^{r,e})$. Especially, $u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^e, \sigma_{-i}^{e+1}) > 0$ indicates there is a strictly policy space expanding at $e + 1$, i.e., $\pi_{-i}^{e+1} \in \Pi_{-i}^{r,e+1} \setminus \Pi_{-i}^{r,e}$.*

*Proof.* Note that the proof is non-trivial since $\Pi_{-i}^{r,e} \subset \Pi_{-i}^{r,e+1}$, so we cannot directly derive Theorem B.1 with the Nash equilibrium $(\pi_i^e, \sigma_{-i}^e)$. Instead, we should combine equilibriums of epoch $e$ and $e + 1$. Additionally, the equilibrium considered is mixed-strategy Nash equilibrium. Considering the property of Nash equilibrium, $\forall \pi_i \in \Delta(\Pi_i), u_i(\pi_i, \sigma_{-i}^{e+1}) \leq u_i(\pi_i^{e+1}, \sigma_{-i}^{e+1})$, then we have

$$u_i(\pi_i^e, \sigma_{-i}^{e+1}) \leq u_i(\pi_i^{e+1}, \sigma_{-i}^{e+1}). \tag{4}$$

Analogously, $\forall \sigma_{-i} \in \Delta_{\Pi_{-i}^r}^{e+1}, u_i(\pi_i^{e+1}, \sigma_{-i}^{e+1}) \leq u_i(\pi_i^{e+1}, \sigma_{-i})$, then we have

$$u_i(\pi_i^{e+1}, \sigma_{-i}^{e+1}) \leq u_i(\pi_i^{e+1}, \sigma_{-i}^e). \tag{5}$$

Combining Eq. (4) and (5), we can derive

$$u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^e, \sigma_{-i}^{e+1}) \tag{6}$$
$$\geq u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^{e+1}, \sigma_{-i}^{e+1})$$
$$\geq u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^{e+1}, \sigma_{-i}^e).$$

Since $(\pi_i^e, \sigma_{-i}^e)$ is a Nash equilibrium, then there is $u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^{e+1}, \sigma_{-i}^e) \geq 0$. Thus $u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^e, \sigma_{-i}^{e+1}) \geq 0$. Then we give a proof of strict expanding as Corollary B.2. $\square$

**Corollary B.2** (*Strict Expanding*). *For equilibrium $(\pi_i^e, \sigma_{-i}^e)$ and $(\pi_i^{e+1}, \sigma_{-i}^{e+1})$, there is a strict expanding when $u_i(\pi_i^e, \sigma_{-i}^e) > u_i(\pi_i^e, \sigma_{-i}^{e+1})$ and $\Pi_{-i}^{e,r} \subset \Pi_{-i}^{e+1,r}$.*

*Proof.* Suppose $\sigma_{-i}^{e+1} \in \Delta_{\Pi_{-i}^r}^e$ when $u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^e, \sigma_{-i}^{e+1}) > 0$ holds. As $(\pi_i^{e+1}, \sigma_{-i}^{e+1})$ is also an equilibrium of $(\Pi_i, \Pi_{-i}^{e,r})$, so there is $u_i(\pi_i^e, \sigma_{-i}^e) = u_i(\pi_i^{e+1}, \sigma_{-i}^{e+1})$. Thus we have

$$u_i(\pi_i^e, \sigma_{-i}^e) - u_i(\pi_i^e, \sigma_{-i}^{e+1}) \geq u_i(\pi^e, \sigma_{-i}^e) - u_i(\pi^{e+1}, \sigma_{-i}^{e+1}) = 0, \tag{7}$$

which violates the inequality. So there must be $\sigma_{-i}^{e+1} \in \Delta_{-i}^{e+1} \setminus \Delta_{-i}^e$. $\square$

**Discussion of Exploration Efficiency**

We further discuss the exploration efficiency from the theoretical perspective as follows.

**Proposition B.3.** *EPSRO has higher exploration efficiency than PSRO.*

*Proof.* Before the proof, we note that $\sigma_i \in \Delta(\Pi_i^r)$ and $\pi_i \in \Delta(\Pi_i)$ for each player $i \in \{1, 2\}$. The NASHCONV (NC) of PSRO is

$$\mathrm{NC}^{psro} = \max_{\pi_i \in \Delta(\Pi_i)} u_i(\pi_i, \sigma_{-i}) + \max_{\pi_{-i} \in \Delta(\Pi_{-i})} u_{-i}(\sigma_i, \pi_{-i})$$
$$= \max_{\pi_i \in \Delta(\Pi_i)} u_i(\pi_i, \sigma_{-i}) - \max_{\pi_{-i} \in \Delta(\Pi_{-i})} u_i(\sigma_i, \pi_{-i})$$
$$= u_i(\pi_i^*, \sigma_{-i}) - u_i(\sigma_i, \pi_{-i}^*),$$

and $u_i(\pi_i^*, \sigma_{-i}) \geq u_i(\sigma_i, \sigma_{-i}) \geq u_i(\sigma_i, \pi_{-i}^*)$.

At each epoch, EPSRO directly learns a tuple of strategies as an equilibrium for each player, i.e., $(\pi_i^*, \sigma_{-i}^*)$ and $(\sigma_i^*, \pi_{-i}^*)$ are two NEs of $(\Pi_i, \Pi_{-i}^r)$, $(\Pi_i^r, \Pi_{-i})$, respectively. Thus, we can apply the same rule as PSRO to derive EPSRO's NASHCONV which satisfies the following inequality:

$$\mathrm{NC}^{epsro} = u_i(\pi_i^*, \sigma_{-i}^*) - u_i(\sigma_i^*, \pi_{-i}^*)$$
$$\leq u_i(\pi_i^*, \sigma_{-i}) - u_i(\sigma_i^*, \pi_{-i}^*)$$
$$\leq u_i(\pi_i^*, \sigma_{-i}) - u_i(\sigma_i, \pi_{-i}^*) = \mathrm{NC}^{psro}.$$

Though $\mathrm{NC}^{epsro} \leq \mathrm{NC}^{psro}$, the equation only holds when $(\pi_i^*, \sigma_{-i})$ and $(\sigma_i, \pi_{-i}^*)$ are equilibrium because it requires $\Pi_i^r = \Pi_i$ for each player. Thus, $\mathrm{NC}^{epsro} < \mathrm{NC}^{psro}$ at each epoch before convergence, i.e. EPSRO has higher exploration than PSRO. $\square$

**Modeling The Learning as MWU**

In this paper, we argue that the updates of meta-strategies (Algorithm 5) follow the rule of Multiplicative Weights Update (MWU) algorithms. Before the explanation of equivalence, we introduce the definition of MWU as follows.

**Definition B.4** (*Multiplicative Weights Update [12]*). Given a game with utility matrix $U$ for the row player, let $c_1, c_2, \ldots$ be a sequence of mixed strategies played by the column player. The row player is said to follow MWU if $\pi_{t+1}$ is updated as follows:

$$\pi_{t+1}(i) = \pi_t(i) \frac{\exp\left(\mu_t {a^i}^\top U c_t\right)}{\sum_{i=1}^n \pi_t(i) \exp\left(\mu_t {a^i}^\top U c_t\right)}, \forall i \in [n] \tag{8}$$

where $\mu_t > 0$ is a parameter, $\pi_0 = [1/n, \ldots, 1/n]$ and $n$ is the number of pure strategies (a.k.a. experts).

Then we prove the equivalence in the following Lemma.

**Lemma B.5.** *Let $\beta = [\beta_1, \ldots, \beta_n]$ be the weights vector of the meta strategy played by player $-i$. The player is said to follow MWU if $\sigma_{-i,\beta'}$ is updated as follows*

$$\sigma_{-i,\beta'}(j) = \frac{\exp\left(\beta_j - g_j\right)}{\sum_{j=1}^n \exp\left(\beta_j - g_j\right)}, \tag{9}$$

*where $g_j$ is the estimated gradients to item $\beta_j$.*

*Proof.* Since the learning for each opponent $-i$ follows the same rule, we mitigate the index $-i$ here.

$$\begin{aligned}
\sigma_{\beta'}(j) &= \frac{\exp\left(\beta_j - g_j\right)}{\sum_{i=1}^n \exp\left(\beta_i - g_i\right)} = \exp\beta_j \frac{\exp-g_j}{\sum_{i=1}^n \exp\beta_i \cdot \exp-g_i} \\
&= \frac{\exp\beta_j}{\sum_{i=1}^n \exp\beta_i} \cdot \frac{\exp-g_j}{\sum_{i=1}^n \frac{\exp\beta_i}{\sum_{k=1}^n \exp\beta_k} \exp-g_i} \\
&= \sigma_\beta(j) \cdot \frac{\exp-g_j}{\sum_{i=1}^n \sigma_\beta(i) \exp-g_i}
\end{aligned}$$

where $g_j = \eta \nabla_{\beta_j} u_{-i}(\pi_i, \sigma_{-i,\beta})$, $\eta$ is the learning rate. $\square$

**A Proof of Theorem 3.7**

As introduced in the main content, we consider a warm-start technique to reduce the re-training cost of meta-strategies and save training steps. To find a feasible meta-strategy $\bar{\sigma}'_{-i}$ for warm-starting, we need to ensure that $\bar{\sigma}'_{-i}$ follows two conditions proposed by Noam Brown et al. (2016) [6], i.e. (1) cannot violate the bound of regret of the **latest epoch** $e + 1$ and (2) cannot violate the Nash equilibrium of the **latest epoch** $e + 1$. In our paper, we choose the average strategies $\bar{\sigma}'_{-i}$ that satisfies $u_i^{e+1}(\bar{\pi}_i, \bar{\sigma}'_{-i}) = u_i^e(\bar{\pi}_i, \bar{\sigma}_{-i})$. Note that the utility function of epoch $e + 1$ is different from the one of epoch $e$ since the dimension of opponent strategy is changed with a new introduced best response $\pi_{-i}^{e+1}$, i.e. $u_i^e(\bar{\pi}_i, \cdot) \in \mathbb{R}^k$ while $u_i^{e+1}(\bar{\pi}_i, \cdot) \in \mathbb{R}^{k+1}$.

**Theorem B.6.** *Suppose a substitute policy of $\bar{\sigma}_{-i} \in \Delta_{\Pi^r_{-i}}^e$ is $\bar{\sigma}'_{-i} \in \Delta_{\Pi^r_{-i}}^{e+1}$, and it satisfies $u_i^{e+1}(\bar{\pi}_i, \bar{\sigma}'_{-i}) = u_i^e(\bar{\pi}_i, \bar{\sigma}_{-i})$, we have*

$$\max_{\sigma_{-i} \in \Delta_{\Pi^r_{-i}}^{e+1}} \sum_{t=1}^T \left( u_{-i}^{e+1}(\pi_i^t, \sigma_{-i}) - u_{-i}^{e+1}(\pi_i^t, \bar{\sigma}'_{-i}) \right) \le \epsilon_{-i}, \tag{10}$$

*where $\epsilon_{-i}$ is the regret bound of epoch $e + 1$, $\pi_i^t \in \Delta(\Pi_i)$.*

*Proof.* Noam Brown et al. have proved in [6] that (1) *the growth of substitute regret has the same bound as the growth rate of normal regret ([6, Lemma. 2])* and (2) *we can use a substitute strategy*

15

*regret to prove convergence to a Nash Equilibrium just as we could use normal regret ([6, Lemma. 3]).* Thus, if we use a substitute strategy $(\bar{\pi}_i, \bar{\sigma}'_{-i})$ that warm-start to $T$ iterations as Eq. 10, and play more $T' \geq 0$ iterations according to Algorithm 5, then we have:

$$\max_{\sigma_{-i} \in \Delta_{\Pi^r_{-i}}^{e+1}} \left( T \left( u_{-i}^{e+1}(\bar{\pi}_i, \sigma_{-i}) - u_{-i}^{e+1}(\bar{\pi}_i, \bar{\sigma}'_{-i}) \right) + \sum_{t'=1}^{T'} \left( u_{-i}^{e+1}(\pi_i^{t'}, \sigma_{-i}) - u_{-i}^{e+1}(\pi_i^{t'}, \sigma_{-i}^{t'}) \right) \right) \leq \epsilon_{-i}.$$
(11)

As the summation of regret from $T + 1$ to $T + T'$ is non-negative, then we know Eq. 10 still holds. Since $\max_{\sigma_{-i}} u_{-i}^{e+1}(\bar{\pi}_i, \sigma_{-i}) \geq u_{-i}^{e+1}(\bar{\pi}_i, \bar{\sigma}'_{-i})$ and $\max_{\pi_i} u_i^{e+1}(\pi_i, \bar{\sigma}_{-i}) \geq u_i^{e+1}(\bar{\pi}_i, \bar{\sigma}_{-i})$, so the substitute strategies $(\bar{\pi}_i, \bar{\sigma}'_{-i})$ satisfies the convergence of $(\epsilon_i, \epsilon_{-i})$-Nash equilibrium of epoch $e + 1$ (See Theorem 3.6 or G.1). Then $(\bar{\pi}_i, \bar{\sigma}'_{-i})$ are feasible warm-start strategies. $\square$

## A Proof of Lemma 3.8

As for the exact bound of regret, we will give a proof in Theorem B.6. Note that the utility function of epoch $e + 1$ is different from the one of epoch $e$ since the dimension of opponent strategy is changed with a new introduced best response $\pi_{-i}^{e+1}$, i.e. $u_{-i}^e(\bar{\pi}_i, \cdot) \in \mathbb{R}^k$ while $u_{-i}^{e+1}(\bar{\pi}_i, \cdot) \in \mathbb{R}^{k+1}$. Thus, before the optimization of $\beta$, we need to slightly run a simulation to estimate the item $u_{-i}(\bar{\pi}_i, \pi_{-i}^{e+1})$ in $u_{-i}^{e+1}$.

**Lemma B.7.** *Let* $k = |\Pi_{-i}^{r,e+1}|$, $\bar{\sigma}'_{-i}$ *be parameterized by* $\beta_{-i} = [\beta_{-i,1}, \beta_{-i,2}, \ldots, \beta_{-i,k}]$, $\bar{\sigma}'_{-i}(k)$ *the $k$-th item of* $\bar{\sigma}'_{-i}$, $x = [\bar{\sigma}'_{-i}(1), \ldots, \bar{\sigma}'_{-i}(k-1)]^\top$, $\bar{l}_{-i}^e$ *is $-i$'s average loss vector to* $\bar{\pi}_i$ *at epoch $e$. Then a feasible initial of* $\beta_{-i}^{e+1}$ *could satisfy*

$$\beta_{-i}^{e+1} = \arg\min_{\beta_{-i}} \| (x - \bar{\sigma}_{-i})^\top \bar{l}_{-i}^e - \bar{\sigma}'_{-i}(k) u_{-i}(\bar{\pi}_i, \pi_{-i}^{e+1}) \|_2 .$$
(12)

*Proof.* Considering the error between $\xi = u_{-i}^e(\bar{\pi}_i, \bar{\sigma}_{-i}) - u_{-i}^{e+1}(\bar{\pi}_i, \bar{\sigma}'_{-i})$. It can be rewritten as

$$\xi = -\bar{\sigma}_{-i}^\top \bar{l}_{-i}^e - \left( -x^\top \bar{l}_{-i}^e + \bar{\sigma}'_{-i}(k) u_{-i}(\bar{\pi}_i, \pi_{-i}^{e+1}) \right)$$

$$= (x - \bar{\sigma}_{-i})^\top \bar{l}_{-i}^e - \bar{\sigma}'_{-i}(k) u_{-i}(\bar{\pi}_i, \pi_{-i}^{e+1}).$$

Thus, we can optimize $\beta_{-i}$ by minimizing the error $\xi$ as $\beta_{-i}^{e+1} = \arg\min_{\beta_{-i}} \| \xi \|_2$ $\square$

We compute the initial of $\sigma_{-i}^{e+1}$ as Lemma B.7. Note that there is another solution as $\bar{\sigma}'_{-i}(k) = 0, \bar{\sigma}_{-i} = x$. In that case, the newly introduced policy $\pi_{-i}^{e+1}$ will have no chance to be sampled, causing biased or even wrong solution. Thus, we consider adding a regularization to Eq. 12 that maximizes the entropy of $\bar{\sigma}'_{-i}$, i.e. $H(\bar{\sigma}'_{-i}) = -\bar{\sigma}'^\top_{-i} \log \bar{\sigma}'_{-i}$ and $\beta_{-i}^{e+1} = \arg\min_{\beta_{-i}} \| \xi \|_2 - \lambda H(\bar{\sigma}'_{-i})$, where $\lambda > 0$.

## A Proof of Theorem 3.9

**Theorem B.8** (*Regret Bound of EPSRO*). *Let* $l_1, l_2, \ldots, l_T$ *be a sequence of loss vectors player by an adversary, and* $\langle \cdot, \cdot \rangle$ *be the dot product, then Algorithm 5 is a no-regret algorithm with*

$$\frac{1}{T} \left( \sum_{t=1}^{T} \langle \sigma_t, l_t \rangle - \min_{\sigma \in \Delta(\Pi_t^r)} \sum_{t=1}^{T} \langle \sigma, l_t \rangle \right) \leq \frac{\sqrt{\log [(k+1)k/2]}}{\sqrt{2T}}, \text{ where $k$ is the size of $\Pi^r$.}$$

*Proof.* Note that the loss vector is the opposite of utility at each iteration, i.e., $l_t = -u_t$. The proof follows the approach of Le Cong Dinh et al. (2021) [20]. Denote $|T_1|, \ldots, |T_k|$ be the learning horizon of each epoch. During each $|T_j|$, the restricted policy set is denoted as unchanged $\Pi^{r,i}$. In the case of finite $k$, we have:

$$\sum_{i=1}^{k} |T_i| = T.$$
(13)

For each training epoch $i$, following the regret bound of MWU, we have:

$$\sum_{t=|T^i|+1}^{|T^{i+1}|} \langle \sigma_t, l_t \rangle - \min_{\sigma \in \Delta(\Pi_i^r)} \sum_{t=|T^i|+1}^{|T^{i+1}|} \langle \sigma, l_t \rangle \leq \sqrt{\frac{|T_i| \log |\Pi^{r,i}|}{2}}, \text{ where } |T^i| = \sum_{j=1}^{i} |T_j|. \quad (14)$$

Considering all time horizon, for $i = 1, \ldots, k$, we have:

$$\sum_{i=1}^{k} \sqrt{\frac{|T_i| \log |\Pi^{r,i}|}{2}} \geq \sum_{t=1}^{T} \langle \sigma_t, l_t \rangle - \sum_{i=1}^{k} \min_{\sigma \in \Delta(\Pi^{r,i})} \sum_{t=|T^i|+1}^{|T^{i+1}|} \langle \sigma, l_t \rangle \quad (15)$$

$$\geq \sum_{t=1}^{T} \langle \sigma_t, l_t \rangle - \min_{\sigma \in \Delta(\Pi^r)} \sum_{i=1}^{k} \sum_{t=|T^i|+1}^{|T^{i+1}|} \langle \sigma, l_t \rangle \quad (16)$$

$$= \sum_{t=1}^{T} \langle \sigma_t, l_t \rangle - \min_{\sigma \in \Delta(\Pi^r)} \sum_{t=1}^{T} \langle \sigma, l_t \rangle. \quad (17)$$

Note that $i = |\Pi^{r,i}|$, then $\sum_{i=1}^{k} \sqrt{|T_i| \log |\Pi^{r,i}|} \leq \sum_{i=1}^{k} \sqrt{T \log(i)}$. As the sum of logarithms is concave, so we have $\sum_{i=1}^{k} \sqrt{T \log(i)} \leq \sqrt{T \log \left( \sum_{i=1}^{k} i \right)} = \sqrt{T \log [(k+1)k/2]}$ then we have

$$\frac{\sqrt{\log [(k+1)k/2]}}{\sqrt{2T}} \geq \frac{1}{T} \left( \sum_{t=1}^{T} \langle \sigma_t, l_t \rangle - \min_{\sigma \in \Delta(\Pi^r)} \sum_{t=1}^{T} \langle \sigma, l_t \rangle \right).$$

$\square$

In addition, we provide a comparison of regret bound on existing solvers in Appendix G.2.

## A Proof of Theorem 3.10

Before the proof of convergence rate of EPSRO, we want to clarify that the reinforcement learning procedure for the best response $\pi_i$ should be an equivalence of MWU. Fortunately, many reinforcement learning algorithms can be treated like that as their policy function is modeled as a Boltzman distribution [7]. Thus, EPSRO can implement such a procedure with algorithms like Soft Actor Critic [15], Soft Q-learning [14], DQN with Boltzman sampling [38] and others. With this condition, we can apply the Theorem B.6 to $\pi_i$ and then derive the convergence rate as follows.

**Theorem B.9** (*Convergence Rate of EPSRO*). *Let $k, N$ denote the size of restricted policy sets $\Pi_{-i}^r$ and $\Pi_i$. Then the learning of Algorithm 3 will converge to the Nash equilibrium with the rate:*

$$\epsilon_T = \sqrt{\frac{\log [(k+1)k/2]}{2T}} + \sqrt{\frac{\log [(N+1)N/2]}{2T}}.$$

*Proof.* Using the regret bound of Theorem B.8 we have:

$$\max_{\pi_i} \frac{1}{T} \left( \sum_{t=1}^{T} u_i(\pi_i, \sigma_{-i}^t) - u_i(\pi_i^t, \sigma_{-i}^t) \right) \leq \epsilon_i,$$

$$\max_{\sigma_{-i}} \frac{1}{T} \left( \sum_{t=1}^{T} u_{-i}(\pi_i^t, \sigma_{-i}) - u_{-i}(\pi_i^t, \sigma_{-i}^t) \right) \leq \epsilon_{-i},$$

where $\epsilon_i = \sqrt{\frac{\log [(N+1)N/2]}{2T}}$ and $\epsilon_{-i} = \sqrt{\frac{\log [(k+1)k/2]}{2T}}$. From the above inequalities and $u_{-i}(\cdot, \cdot) = -u_i(\cdot, \cdot)$ in zero-sum games, we can derive that

$$u_i(\bar{\pi}_i, \bar{\sigma}_i) \geq \min_{\sigma_{-i} \in \Delta(\Pi_{-i}^r)} u_i(\bar{\pi}_i, \sigma_{-i}) \geq \frac{1}{T} \sum_{t=1}^{T} u_i(\pi_t, \sigma_t) - \epsilon_{-i}$$

$$\geq \max_{\pi_i} u_i(\pi_i, \bar{\sigma}_{-i}) - \epsilon_i - \epsilon_{-i}.$$

17

By symmetry, we have

$$u_i(\bar{\pi}_i, \bar{\sigma}_i) \leq \max_{\pi_i} u_i(\pi_i, \bar{\sigma}_{-i}) \leq \frac{1}{T} \sum_{t=1}^{T} u_i(\pi_i^t, \sigma_{-i}^t) + \epsilon_i$$
$$\leq \min_{\sigma_{-i}} u_i(\bar{\pi}_i, \sigma_{-i}) + \epsilon_{-i} + \epsilon_i.$$

Thus, with $\epsilon_T = \epsilon_i + \epsilon_{-i}$, we have

$$\max_{\pi_i} u_i(\pi_i, \bar{\sigma}_{-i}) - \epsilon_T \leq u_i(\bar{\pi}_i, \bar{\sigma}_{-i}) \leq \min_{\sigma_{-i}} u_i(\bar{\pi}_i, \sigma_{-i}) + \epsilon_T.$$

By definition, we have $(\bar{\pi}_i, \bar{\sigma}_{-i})$ is $\epsilon_T$-Nash equilibrium. $\square$

## C  Implementation and Hyper-parameter Selection

**Learning Meta-strategies and Best-response.**  We introduce the algorithm for learning meta-strategies in Algorithm 5. We train best response policies with off-policy reinforcement learning algorithm, DQN. For all involved PSRO baselines in our paper, the selected implementation for best response learning is DQN.

**Parameter Selection.**  We keep the consistency on the implementation of policy support in each PSRO-based method in this paper. Expressly, the network is set to 4 Dense layers, 256 units each. The learning rate for reinforcement learning policy is set to 0.01. For the hyper-parameters of meta-strategy, the window size $L$ is set to 100 episodes, the learning rate is 0.01 for Kuhn Poker, and 0.005 for other environments. The number of parallel workers is set to 4 for EPSRO and P-PSRO in all experiments.

## D  Baselines

**Self-Play.**  Self-play is an open-ended learning algorithm for multi-agent reinforcement learning [17]. In the training process, self-play generates a sequence of policies and keeps training policies against the newest opponents. This algorithm outperforms in some classic games, such as Go and Chess. However, self-play fails in nontransitive games.

**Policy Space Response Oracles (PSRO).**  PSRO algorithm is well described above in previous sections. It provides an iterative solution to solve the approximation of Nash equilibrium for large games [19]. PSRO iteratively trains new policies against a meta-strategy of opponent population and expends policy populations with the current well-trained policy.

**Rectified PSRO (PSRO-rN).**  PSRO-rN is a variant of PSRO that aims to solve non-transitive zero-sum games [2], such as rock-paper-scissors. This algorithm involves rectified Nash response to construct adaptive sequences of objectives for non-transitive games. Policies in PSRO-rN only train against others that they already beat.

**Mixed Oracles.**  Mixed Oracles is another variant of PSRO that aims to improve computational efficiency by reducing training costs [32]. At each iteration, it utilizes knowledge of former iterations, thus only needing to train current policies against the newest opponent.

**Pipeline-PSRO (P-PSRO).**  To further accelerate the training process of PSRO, P-PSRO is proposed to parallelize the training process [24]. Compared to other parallel algorithms, such as DHC, which fail to converge in some cases, P-PSRO maintains a parallel pipeline of learning workers with convergence guarantees.

## E  Environment Details

We introduce more details about the random symmetric games and mulit-agent gathering environments here.

**Random Symmetric Games.** McAleer et al. [24] introduce the games to investigate the performance of PSRO-based methods in high-dimensional symmetric games (SymGame). In this experiment, we generated random symmetric zero-sum matrices with different dimension $n$. For a given matrix, elements in the upper triangle are distributed uniformly: $\forall i < j \le n, a_{i,j} \sim \text{UNIFORM}(-1, 1)$ and for the lower triangle, the elements are set to be the negative of its diagonal counterpart: $\forall j < i \le n, a_{i,j} = -a_{j,i}$. The diagonal elements are equal to zero: $a_i, i = 0$. The matrix defines the utility of two pure strategies to the row player. A strategy $\pi \in \Delta^n$ is a distribution over the $n$ pure strategies of the game given by the rows (or equivalently, columns) of the matrix.

**Multi-agent Gathering.** We introduce two multi-agent gathering environments in this paper, the *Gathering Small* and the *Gathering Open*. For each environment, the agent number is set to 2, and the difference between them is that the *Gathering Open* has a much bigger map than *Gathering Small*. Therefore, the agents need to explore a higher dimensional state space in the *Gathering Open* than the smaller one.

# F   Additional Experimental Results

## F.1   Non-transitive Mixture Game

We test six PSRO-based algorithms in the *non-transitive mixture game* to investigate the exploration efficiency. We also introduce the results of NASHCONV to compare the performance. We find that EPSRO outperforms all other algorithms in this game, and performs the highest exploration efficiency.
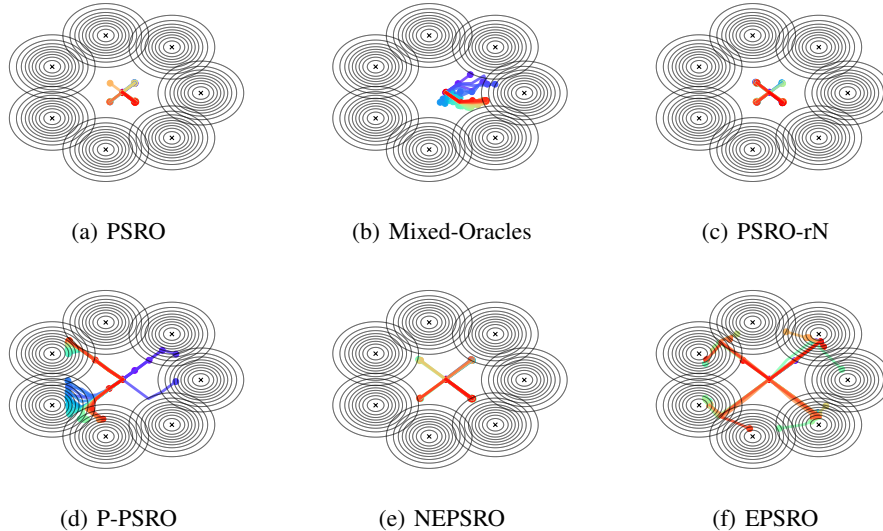


Figure 6: Exploration trajectories on *Non-transitive Mixture Games*. The more trajectories close to the centers of Gaussian, the higher the exploration efficiency of the algorithm. EPSRO outperform all selected baselines since it explored all centers.

## F.2   Random Symmetric Matrix Game

We compare the NASHCONV over iteration of EPSRO with PSRO, P-PSRO, Rectified PSRO, Self-Play, Mixed-Oracles and naive EPSRO without pipeline training (NEPSRO). We run 5 experiments for each set of dimension. The dimensions of size including 15, 30, 45, 60 and 120. The learning rates is set to 0.5, and 4 parallel threads for parallel algorithms. We find that EPSRO performs better than all other algorithms in every dimension setting.
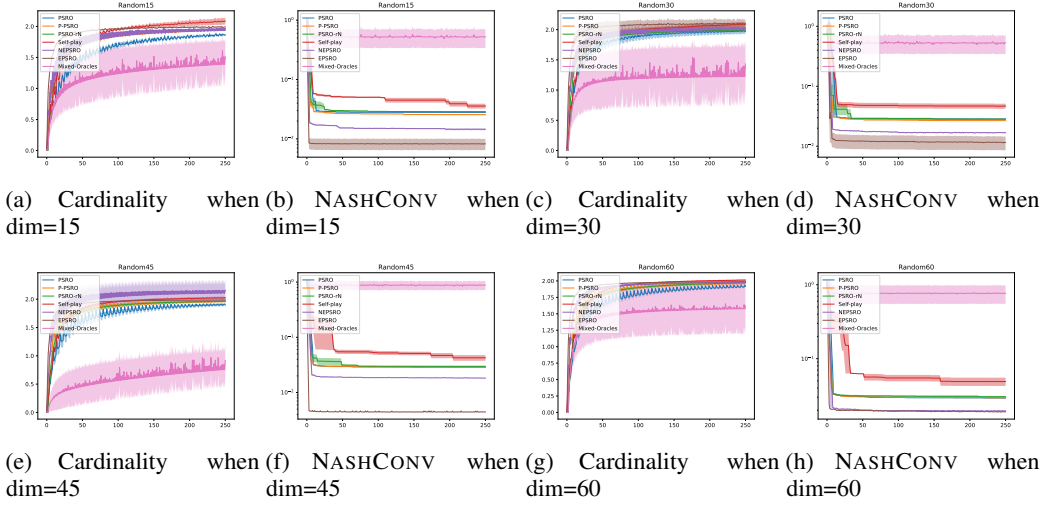
(a) Cardinality when dim=15
(b) NASHCONV when dim=15
(c) Cardinality when dim=30
(d) NASHCONV when dim=30

(e) Cardinality when dim=45
(f) NASHCONV when dim=45
(g) Cardinality when dim=60
(h) NASHCONV when dim=60

Figure 7: Comparison on NASHCONV and cardinality for random symmetric matrix games with different dimensions.

## F.3 Multi-agent Gathering

Multi-agent Gathering environments (MAG) have complex state space than the other games, so that it is highly computation expensive to traverse the game tree to compute the NASHCONV. Instead, we evaluate all algorithms with a fixed policy set. In our experiments, the fixed policy set is generated with PSRO, i.e. $\Pi^{\text{PSRO}}$. We list the pseudo-code for evaluation in Algorithm 7.

---

**Algorithm 7:** EMPIRICAL EVALUATION ON MAG

**Input:** a policy set $\Pi^{\text{TEST}}$ of evaluated algorithm; $\Pi^{\text{PSRO}}$; an empty matrix $M \in \mathbb{R}^{|\Pi^{\text{TEST}}| \times |\Pi^{\text{PSRO}}|}$

1 **for** *each policy* $\pi_i^{\text{TEST}}$ *in* $\Pi^{\text{TEST}}$ **do**

2      **for** *each policy* $\pi_j^{\text{PSRO}}$ *in* $\Pi^{\text{PSRO}}$ **do**

3          Run 50 episodes to evaluate $M_{i,j} = u_i(\pi_i^{\text{TEST}}, \pi_j^{\text{PSRO}})$

4      Compute score of $\sigma_{1:i}^{\text{TEST}}$ as $\text{SCORE}(\sigma_{1:i}^{\text{TEST}}) = \sigma_{1:i}^{\text{TEST}} M_{1:i} \left[ \sigma^{\text{PSRO}} \right]^T$

**Output:** a list of score $\text{SCORE}(\Pi^{\text{TEST}}) = \{ \text{SCORE}(\sigma_{1:i}^{\text{TEST}}) \mid i = 1, \ldots, |\Pi^{\text{TEST}}| \}$ for $\Pi^{\text{TEST}}$

---

$\sigma_{1:i}^{\text{TEST}}$ in Algorithm 7 indicates a meta-strategy composed of $\pi_1, \ldots, \pi_i$, and $M_{1:i}$ is a sub matrix with row 1 to $i$. We present $\text{SCORE}(\Pi^{\text{TEST}})$ of each algorithm in Figure 5. The curve of Self-Play is not included because we keep only two policies (one for the opponent, another for training) in our implementation.

In the training stage, we set the number of simulations for each joint policy as 100, and 10000 episodes to optimize each policy. Except for EPSRO, the training time of each algorithm on the Gathering Small is about 24 hours, and 26 hours for the Gathering Open, while the training time for EPSRO is about 8 hours.

# G Additional

## G.1 Convergence Proof of Average Substitute Strategies

**Theorem G.1.** *([6, Theorem 1]). In a two-player zero-sum game, if $\frac{R_i^T}{T} \leq \epsilon_i$ for both player $i \in \{1, 2\}$, then $\bar{\sigma}$ is a $(\epsilon_1 + \epsilon_2)$-equilibrium.*

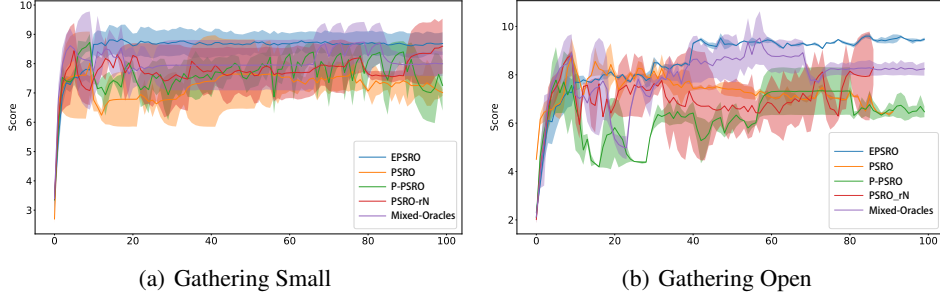(a) Gathering Small              (b) Gathering Open

Figure 8: The score of algorithms on two multi-agent gathering environments. The horizon axis indicates the number of training iterations. As reported in this Figure, EPSRO performs better than other algorithms.

*Proof.* Follow the proof approach of Waugh et al. (2009). From (5), we have that

$$\max_{\sigma \in \Delta} \frac{1}{T} \left( \sum_{t=1}^{T} u_i(\sigma'_i, \sigma^t_{-i}) - u_i(\sigma^t_i, \sigma^t_{-i}) \right) \le \epsilon_i. \tag{18}$$

Since $\sigma'_i$ is the same on every iteration, this becomes

$$\max_{\sigma' \in \Delta} u_i(\sigma'_i, \bar{\sigma}^T_{-i}) - \frac{1}{T} \sum_{t=1}^{T} u_i(\sigma^t_i, \sigma^t_{-i}) \le \epsilon_i. \tag{19}$$

Since $u_i(\sigma) = u_2(\sigma)$, if we sum Equation 19 for both players

$$\max_{\sigma'_1 \in \Delta} u_1(\sigma'_1, \bar{\sigma}^T_2) + \max_{\sigma'_2 \in \Delta} u_2(\bar{\sigma}^T_1, \sigma'_2) \le \epsilon_1 + \epsilon_2, \tag{20}$$

$$\max_{\sigma'_1 \in \Delta} u_1(\sigma'_1, \bar{\sigma}^T_2 - \min_{\sigma'_2 \in \Delta} u_1(\bar{\sigma}^T_1, \sigma'_2) \le \epsilon_1 + \epsilon_2. \tag{21}$$

Since $u_1(\bar{\sigma}^T_1, \bar{\sigma}^T_2) \ge \min_{\sigma'_2 \in \Delta} u_1(\bar{\sigma}^T_1, \sigma'_2)$, so we have $\max_{\sigma'_1 \in \Delta} u_1(\sigma'_1, \bar{\sigma}^T_2) - u_1(\bar{\sigma}^T_1, \bar{\sigma}^T_2) \le \epsilon_1 + \epsilon_2$. By symmetry, this is also true for player 2. Therefore, $\langle \bar{\sigma}^T_1, \bar{\sigma}^T_2 \rangle$ is a $(\epsilon_1 + \epsilon_2)$-equilibrium. □

When warm starting, we can calculate this value because we set $\bar{\sigma}^T = \sigma$. However we cannot calculate $\sum_{t=1}^{T} u_i(\sigma^t)$ because we did not define individual strategies played on each iteration. Fortunately, it turns out we can substitute another value we refer to as $Tu_i(\bar{\sigma})$, chosen from a range of acceptable options. To see this we first observe that the value of $\sum_{t=1}^{T} u_i(\sigma^t)$ is not relevant to the proof of Theorem G.1. Specifically, in Eq. 20, we see it cancels out. Thus, if we choose $(\bar{\pi}_o, \bar{\sigma}_{-i})$ such that satisfies it. Since $\max_\pi u_i(\pi, \bar{\sigma}^T_{-i}) \ge u_i(\bar{\pi}_i, \bar{\sigma}_{-i})$ and $\max_{\sigma_{-i}} u_{-i}(\bar{\pi}_i, \sigma_{-i}) \ge u_{-i}(\bar{\pi}_i, \bar{\sigma}_{-i})$. Thus $(\bar{\pi}_i, \bar{\sigma}_{-i})$ is a feasible warm-starting strategy.

## G.2   Comparison of Regret Bound

We compare EPSRO with existing solvers in the Table. Properties considered for comparison including (1) time complexity or regret bound; (2) whether we need to do retraining; (3) whether we need to know the full game.

| Method | Time Complexity ($\tilde{\mathcal{O}}$) / Regret Bound ($\mathcal{O}$) | No Need to Re-train | No Need to Know the Full Game |
|---|---|---|---|
| Linear Programming [34] | $\tilde{\mathcal{O}}(n \exp(-T/n^2 .38))$ | $\times$ | $\times$ |
| Fictitious Play [21] | $\tilde{\mathcal{O}}(T^{-1/(n+m-2)})$ | $\times$ | $\checkmark$ |
| Double Oracle [27] | $\tilde{\mathcal{O}}(n \exp(-T/n^3 .88))$ | $\times$ | $\checkmark$ |
| Multipli. Weight Update[12] | $\mathcal{O}(\sqrt{\log n/T})$ | $\times$ | $\checkmark$ |
| Policy Response Oracles [19] | $\times$ | $\times$ | $\checkmark$ |
| Online Double Oracle [20] | $\mathcal{O}(\sqrt{k \log k/T})$ | $\times$ | $\checkmark$ |
| **EPSRO** | $\mathcal{O}(\sqrt{\log[(k^2+k)/2]/T})$ | $\checkmark$ | $\checkmark$ |