

GENERATIVE EVOLUTIONARY META-SOLVER (GEMS): SCALABLE SURROGATE-FREE MULTI-AGENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Scalable multi-agent reinforcement learning (MARL) remains a central challenge for AI. Existing population-based methods, like Policy-Space Response Oracles, PSRO, require storing explicit policy populations and constructing full payoff matrices, incurring quadratic computation and linear memory costs. We present Generative Evolutionary Meta-Solver (GEMS), a surrogate-free framework that replaces explicit populations with a compact set of latent anchors and a single amortized generator. Instead of exhaustively constructing the payoff matrix, GEMS relies on unbiased Monte Carlo rollouts, multiplicative-weights meta-dynamics, and a model-free empirical-Bernstein UCB oracle to adaptively expand the policy set. Best responses are trained within the generator using an advantage-based trust-region objective, eliminating the need to store and train separate actors. We evaluated GEMS in a variety of Two-player and Multi-Player games such as the Deceptive Messages Game, Kuhn Poker and Multi-Particle environment. We find that GEMS is up to $6\times$ faster, has $1.3\times$ less memory usage than PSRO, while also reaps higher rewards simultaneously. These results demonstrate that GEMS retains the game theoretic guarantees of PSRO, while overcoming its fundamental inefficiencies, hence enabling scalable multi-agent learning in multiple domains.

1 INTRODUCTION

Imagine organizing a round-robin tennis tournament with hundreds of players (Fig. 1). Scheduling every match gives a complete ranking but is clearly inefficient: the number of matches grows quadratically, and the results table quickly becomes unwieldy.

Training AI agents in two-player zero-sum Markov games faces a similar challenge. Population-based methods such as Policy-Space Response Oracles (PSRO) (Lanctot et al. (2017)) maintain a growing set of k policies and explicitly construct a $k \times k$ payoff matrix, where each entry records the expected outcome of a policy against another. A meta-solver then updates the distribution over policies, analogous to computing player rankings from the tournament table.

While conceptually straightforward, classical PSRO suffers from three critical bottlenecks: ❶ *Memory overhead*: storing a new policy for each player leads to linear growth in storage; ❷ *Computation overhead*: filling in the full $k \times k$ payoff matrix quickly becomes infeasible; ❸ *Scalability of new entries*: adding a new policy requires training and storing another separate actor. Prior work has mitigated some of these costs through selective best-response training (Smith et al. (2021)), Double Oracle methods (McAleer et al. (2021); Huang et al. (2022)), meta-game improvements (McAleer et al. (2022b), McAleer et al. (2022a)), and knowledge transfer for new agents (Smith et al. (2023); Lian et al. (2024)). However, these approaches retain the core paradigm of explicit policy sets, leaving scalability fundamentally limited.

We introduce GEMS, a surrogate-free framework that overcomes these limitations while preserving the game-theoretic guarantees of PSRO. Analogous to running a large tournament without scheduling every match, GEMS: ❶ Maintains a compact *anchor set* Z_t of latent codes that implicitly represent active “players”; ❷ Treats the payoff matrix as conceptual and queries it only through unbiased Monte Carlo rollouts, observing a small subset of matches instead of all k^2 ; ❸ Updates the meta-strategy using multiplicative-weights discretization of replicator dynamics, akin to adjusting rankings based on sampled outcomes; ❹ Expands the population with an empirical-

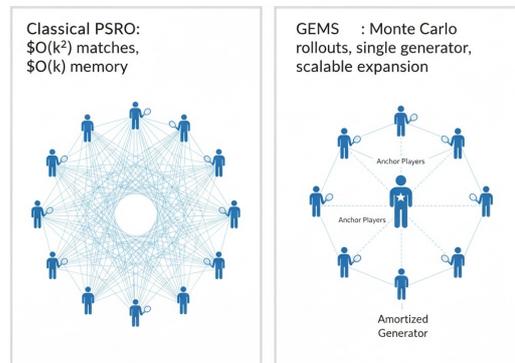


Figure 1: **Tournament analogy for policy populations.** (Left) PSRO: explicit $k \times k$ payoff matrix with all pairwise matchups. (Right) GEMS: compact anchor set of latent policies, with a single generator producing diverse strategies on demand.

Bernstein UCB (EB-UCB) oracle (Maurer & Pontil, 2009), selecting promising new “players” efficiently from a candidate pool; ⑤ Trains a *single* amortized generator G_θ with an Advantage-based Best-Response objective and Trust-region Regularization (ABR-TR), eliminating the need for separate policies.

We evaluated GEMS in a variety of Two-player and Multi-Player games such as the Deceptive Messages Game, Kuhn Poker (Kuhn, 2016) and Multi-Particle environments (Terry et al., 2021). We find that GEMS is up to $6\times$ faster, has $1.3\times$ less memory usage than PSRO, while also reaps higher rewards simultaneously.

Contributions. Relative to classical PSRO, GEMS achieves: ① *Memory efficiency*: replaces $O(k)$ stored players with a single versatile generator; ② *Computation efficiency*: avoids quadratic payoff tables using Monte Carlo estimates, scaling per iteration with sampled matches and candidate pool size; ③ *Scalable new entries*: EB-UCB identifies strong candidates, and ABR-TR integrates them into the generator without adding new actors; ④ *Theoretical guarantees*: unbiased MC meta-gradients, instance-dependent regret for EB-UCB, external regret bounds for multiplicative-weights dynamics, and finite-population exploitability accounting for approximate best responses. Together, these advances transform the exhaustive “tournament bookkeeping” of PSRO into a leaner, scalable framework, closer to how real tournaments operate, where only select matches determine rankings, and strategies adapt efficiently from sparse outcome.

2 RELATED WORK

PSRO and its Variants. The seminal work of Lanctot et al. (2017) introduced PSRO as a general framework for multi-agent reinforcement learning. PSRO iteratively expands a population of policies by training a new best-response policy against a meta-game mixture of the existing population. A core challenge of PSRO is its reliance on a full $k \times k$ payoff matrix, where k is the number of policies in the population. This leads to $O(k^2)$ computation overhead per iteration and $O(k)$ memory overhead to store individual policies, both of which become bottlenecks in large-scale settings (Shao et al., 2024; Zhou et al., 2022). This has motivated a line of work on more efficient variants. For example, Smith et al. (2021) propose training a best-response against a single opponent, reducing computation but potentially limiting diversity. The double oracle (DO) family of algorithms, including XDO (McAleer et al., 2021) and EDO (Huang et al., 2022), focus on finding equilibria more efficiently, often by guaranteeing linear convergence or monotonic exploitability reduction (McAleer et al., 2022b). *However, these methods can still suffer from issues like performance oscillations or slow convergence due to the need for a full set of strategies (Huang et al., 2022). In contrast, GEMS avoids these issues entirely by replacing the explicit payoff matrix with unbiased Monte Carlo rollouts and the discrete set of policies with a single, amortized generator.*

Other work has addressed the computational cost of training new policies from scratch at each iteration, which is a significant bottleneck (McAleer et al., 2020; Li et al., 2023). To address this, FUSION-PSRO (Lian et al., 2024) uses policy fusion to initialize new best responses, while STRATEGIC KNOWLEDGE TRANSFER (Smith et al., 2023) explores transferring knowledge across changing opponent strategies. *Our work tackles this problem differently by using a single amortized generator, which removes the need to store and train separate models for each policy.*

Finally, some papers have focused on improving the meta-game solution or policy diversity. ALPHA-PSRO (Muller et al., 2019) replaces the Nash meta-solver with α -Rank to scale PSRO to general-sum, many-player settings. A-PSRO (Hu et al., 2023) introduces an advantage-based evaluation for strategy selection, providing a unified objective across zero-sum and general-sum games. ANYTIME PSRO (McAleer et al., 2022b) and SELF-PLAY PSRO (SP-PSRO) (McAleer et al., 2022a) aim to improve convergence by adding high-quality policies to the population. Other works introduce new metrics for policy diversity to ensure better approximation of the Nash Equilibrium (Tang et al., 2025; Yao et al., 2023). *While these works aim to improve PSRO, they do not fundamentally change its core structure, which still relies on an explicit population of policies and a payoff matrix.*

Game-Theoretic Methods for Multi-Agent Learning. Beyond the PSRO paradigm, other game-theoretic approaches exist for multi-agent learning. COUNTERFACTUAL REGRET MINIMIZATION (CFR) (Zinkevich et al., 2007) is a well-known method for extensive-form games, but it typically requires explicit state enumeration, making it challenging to scale to large or continuous domains. Some work has explored bridging the gap between PSRO and CFR, as seen in the unified perspective proposed by Wang et al. (2022). Other approaches include methods based on fictitious play, such as FICTITIOUS CROSS-PLAY (FXP) (Xu et al., 2023), which combines self-play and PSRO to find global Nash equilibria. However, many self-play methods lack theoretical guarantees for general-sum games and can struggle with convergence in mixed cooperative-competitive settings (Xu et al., 2023).

Our work contributes to this broader landscape by offering a surrogate-free, amortized framework that is both memory and computationally efficient. While many variants of PSRO have been proposed, they have largely retained the core structure of maintaining an explicit policy population and a large payoff matrix. GEMS breaks

from this paradigm by using an amortized generator and unbiased Monte Carlo rollouts, thereby sidestepping the fundamental scalability issues inherent to classical PSRO.

3 PROPOSED METHOD: GEMS

GEMS circumvents PSRO’s limitations by being **surrogate-free**. Instead of storing k explicit actor models and computing their full payoff matrix, it maintains a single generative model G_θ that maps low-dimensional latent codes to policies. Consequently, the method scales to massive conceptual populations while storing only the generator and a set of latent “anchor” codes.

GEMS unfolds in an iterative loop: it first solves for the equilibrium of the current meta-game using noisy estimates, then expands the game by finding an approximate best response via a bandit-like oracle. The full procedure is detailed below.

While the core exposition focuses on two-player zero-sum games for clarity, the GEMS framework extends naturally to more general settings. This extensibility is demonstrated by benchmarking GEMS on two-player general-sum and, more broadly, **n-player general-sum** games. For the multi-player environments, the implementation leverages the PettingZoo library (Terry et al., 2021). The necessary modifications to the meta-game estimators, per-player oracles, and convergence guarantees are provided in Appendix Part II.

3.1 FORMAL SETUP AND GENERATIVE REPRESENTATION

Consider a two-player zero-sum Markov game. Let $r(\pi_i, \pi_j) \in [0, 1]$ denote the expected return for Player 1 when policy π_i faces π_j . For a finite policy set $A = \{\pi_1, \dots, \pi_k\}$, this induces a payoff matrix $M \in [-1, 1]^{k \times k}$ where $M_{ij} = \mathbb{E}[r(\pi_i, \pi_j)]$. **Crucially, GEMS never explicitly constructs or stores this matrix M .**

At each iteration t , GEMS maintains three core components:

1. An **anchor set** of latent codes $Z_t = \{z_1, \dots, z_{k_t}\} \subset \mathbb{R}^d$, representing policies in a low-dimensional space.
2. A single **generator network** G_θ that maps a latent code z to policy parameters $\varphi = G_\theta(z)$, thereby defining a policy π_φ .
3. A **meta-strategy** $\sigma_t \in \Delta_{k_t-1}$, which is a probability distribution over the current anchor set Z_t , representing the Nash equilibrium of the restricted game $A_t = \{\pi_{G_\theta(z)} : z \in Z_t\}$.

These components define the key quantities for analyzing the meta-game. The vector of expected payoffs for each anchor policy against the meta-strategy σ_t is $v_t = M\sigma_t$, and the expected value of the game at iteration t is $\bar{r}_t = \sigma_t^\top M\sigma_t$. A primary objective is to minimize **exploitability**, the incentive for any player to deviate from σ_t :

$$\text{Exploit}(\sigma_t) = \max_{i \leq k_t} e_i^\top M\sigma_t - \sigma_t^\top M\sigma_t. \quad (1)$$

GEMS iteratively refines σ_t and expands Z_t to drive this exploitability toward zero.

3.2 ESTIMATING THE META-GAME WITHOUT THE MATRIX

Because GEMS does not access the true payoff matrix M , it estimates the values of v_t and \bar{r}_t through simulation. **Adopting a two-time-scale assumption in which the meta-strategy σ_t is held fixed during estimation**, GEMS employs Monte Carlo rollouts to obtain statistically sound estimates.

For each anchor policy $i \in [k_t]$, the algorithm samples n_i opponents $j_s \sim \sigma_t$ and executes m game episodes per pair, yielding returns $Y_{i,s,\ell} \in [0, 1]$. This process yields a simple yet powerful estimator for per-policy performance. To estimate the overall game value, the procedure additionally samples B independent policy pairs $(i_b, j_b) \sim \sigma_t \times \sigma_t$ and averages their outcomes across m episodes each. The resulting estimators are

$$\hat{v}_{t,i} = \frac{1}{n_i m} \sum_{s=1}^{n_i} \sum_{\ell=1}^m Y_{i,s,\ell}, \quad \hat{r}_t = \frac{1}{Bm} \sum_{b=1}^B \sum_{\ell=1}^m Y_{i_b, j_b, \ell}, \quad (i_b, j_b) \sim \sigma_t \times \sigma_t, \quad (2)$$

and constitute the empirical backbone of GEMS.

Lemma 3.1 (Unbiasedness and Empirical-Bernstein Concentration). *With rewards in $[0, 1]$, the estimators are unbiased: $\mathbb{E}[\hat{v}_{t,i}] = (M\sigma_t)_i$ and $\mathbb{E}[\hat{r}_t] = \sigma_t^\top M\sigma_t$. Moreover, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,*

$$|\hat{v}_{t,i} - (M\sigma_t)_i| \leq \sqrt{\frac{2 \widehat{\text{Var}}_{t,i} \ln(2/\delta)}{n_i m}} + \frac{3 \ln(2/\delta)}{n_i m - 1}, \quad |\hat{r}_t - \sigma_t^\top M\sigma_t| = O\left(\sqrt{\frac{\ln(1/\delta)}{Bm}}\right), \quad (3)$$

where $\widehat{\text{Var}}_{t,i}$ is the empirical variance of $\{Y_{i,s,\ell}\}$.

Proof sketch. Unbiasedness follows from the law of total expectation over the sampling of opponents, and the concentration bounds apply the empirical-Bernstein inequality for bounded random variables. \square

3.3 SOLVING THE META-GAME VIA OPTIMISTIC REPLICATOR DYNAMICS

Given the estimated payoffs \hat{v}_t and game value \hat{r}_t , GEMS updates the meta-strategy from σ_t to σ_{t+1} . GEMS adopts the **Optimistic Multiplicative Weights Update (OMWU)** algorithm (Daskalakis & Panageas (2018)), an adaptive discretization of replicator dynamics. Rather than relying solely on the current payoff estimate \hat{v}_t , OMWU incorporates a predictive “hint” about the next payoff. Specifically, it forms the optimistic estimate

$$m_t = 2\hat{v}_t - \hat{v}_{t-1}, \quad \text{with } \hat{v}_0 = \mathbf{0}.$$

The meta-strategy then updates according to

$$\sigma_{t+1}(i) \propto \sigma_t(i) \exp(\eta_t [2\hat{v}_{t,i} - \hat{v}_{t-1,i} - \hat{r}_t]), \quad \eta_t > 0, \quad (4)$$

followed by normalization to ensure $\sigma_{t+1} \in \Delta_{k_t-1}$. This optimistic step yields stronger theoretical guarantees: regret now scales with the cumulative variation of the payoff vectors rather than the iteration horizon T .

Proposition 3.2 (External Regret of OMWU under Unbiased Noise). *Assume payoffs in $[0, 1]$. For any sequence of payoff estimates \hat{v}_t satisfying $\mathbb{E}[\hat{v}_t | \sigma_t] = M\sigma_t$ and a constant step size η , the average external regret obeys*

$$\frac{1}{T} \sum_{t=1}^T \left(\max_i e_i^\top M\sigma_t - \sigma_t^\top M\sigma_t \right) \leq O\left(\frac{1}{T} \sqrt{\ln k_T \sum_{t=1}^T \|v_t - v_{t-1}\|_\infty^2}\right) + \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t - M\sigma_t\|_\infty], \quad (5)$$

where $v_t = M\sigma_t$ are the true expected payoffs. The first term, driven by the variation of the meta-game, leads to faster convergence in slowly changing environments such as those induced by GEMS.

3.4 FINDING NEW STRATEGIES WITH A BANDIT ORACLE

Solving the restricted meta-game alone is insufficient; progress requires introducing new, challenging policies into the population. GEMS casts this search as a **multi-armed bandit problem** (Robbins (1952)). The “arms” comprise a finite pool of candidate latent codes Λ_t , and pulling an arm corresponds to evaluating the policy $G_\theta(z)$ against the current meta-strategy σ_t .

For a latent code $z \in \Lambda_t$, define

$$f_t(z) = \mathbb{E}_{j \sim \sigma_t} [r(\pi_{G_\theta(z)}, \pi_j)]. \quad (6)$$

GEMS estimates this value, $\hat{\mu}_t(z)$, together with its empirical variance $\widehat{\text{Var}}_t(z)$ via Monte Carlo rollouts. To balance exploration and exploitation efficiently, it employs an **empirical-Bernstein Upper Confidence Bound (EB-UCB)** (Maurer & Pontil (2009)), which leverages variance information for tighter bounds. A Jacobian penalty encourages smoothness in the generator’s latent space, aiding optimization. The score assigned to candidate z is

$$\text{UCB}_t^{\text{EB}}(z) = \hat{\mu}_t(z) + \sqrt{\frac{2\widehat{\text{Var}}_t(z) \ln(3/\delta_t)}{n_z m}} + \frac{3 \ln(3/\delta_t)}{n_z m - 1} - \lambda_J \|JG_\theta(z)\|_F^2,$$

where $\delta_t = t^{-2}$ is a decaying confidence parameter. GEMS selects

$$z_t^* = \arg \max_{z \in \Lambda_t} \text{UCB}_t^{\text{EB}}(z),$$

and augments the anchor set: $Z_{t+1} \leftarrow Z_t \cup \{z_t^*\}$.

Theorem 3.3 (Instance-Dependent Oracle Regret). *Assume rewards lie in $[0, 1]$ and the bandit problem admits a unique best arm z^* with sub-optimality gaps $\Delta_z > 0$. Under the two-time-scale assumption (fixed σ_t during selection), the cumulative regret of the oracle satisfies*

$$\sum_{t=1}^T \mathbb{E}[f_t(z^*) - f_t(z_t^*)] = O\left(\sum_{z \neq z^*} \frac{\ln T}{\Delta_z}\right) + \lambda_J \sum_{t=1}^T \mathbb{E}[\|JG_\theta(z_t^*)\|_F^2]. \quad (7)$$

If the Jacobian norm is uniformly bounded, the second term can be controlled by annealing λ_J .

Proof sketch. Standard bandit-analysis arguments apply. With high probability, the EB-UCB provides valid confidence bounds; comparing the chosen arm z_t^* with the optimal arm z^* and summing over time yields the stated instance-dependent regret bound. \square

232 3.5 TRAINING THE GENERATOR WITH AMORTIZED BEST RESPONSE

233
234 Once a promising latent code z_t^* has been identified, GEMS must ensure that the generator G_θ can realize the
235 associated high-performing policy while retaining its ability to produce previously effective policies. This is
236 achieved via an **Amortized Best-Response with a Trust Region (ABR-TR)** objective, inspired by trust-region
237 methods in deep reinforcement learning such as TRPO and PPO (Schulman et al. (2015; 2017)).

238 Rather than training a new network from scratch, GEMS fine-tunes the existing generator G_θ to maximize the
239 performance of a curated set of promising latent codes against the opponent mixture σ_{t+1} . A KL-divergence
240 penalty against a frozen, older generator θ^- serves as a trust region, mitigating catastrophic forgetting:

$$241 \mathcal{L}_{\text{ABR-TR}}(\theta) = \mathbb{E}_{z \sim q_t, j \sim \sigma_{t+1}} \left[\widehat{A}_{\pi_{G_\theta(z)}} - \beta D_{\text{KL}}(\pi_{G_\theta(z)} \parallel \pi_{G_{\theta^-}(z)}) - \lambda_J \|JG_\theta(z)\|_F^2 \right], \quad (8)$$

242 where \widehat{A} denotes a suitable advantage estimator. Taking a few gradient-ascent steps on $\mathcal{L}_{\text{ABR-TR}}$ each iteration
243 amortizes best-response computation while preserving a stable, single generator.

244 3.6 OVERALL EXPLOITABILITY BOUND

245 The preceding components now combine to yield a convergence guarantee for the entire GEMS algorithm. The
246 average exploitability decomposes into four interpretable terms, each tracing back to a distinct error source: (1)
247 the inherent regret of the OMWU meta-strategy solver; (2) noise from Monte-Carlo estimation; (3) sub-optimality
248 of the bandit oracle; and (4) approximation error introduced by amortized generator training.

249 Let ε_{BR} denote the gap between a true best response and the policy produced by the generator for the selected
250 latent code z_t^* :

$$251 \varepsilon_{\text{BR}} = \sup_t \mathbb{E}_{j \sim \sigma_t} \left[r(\pi_{\text{BR}}(j), \pi_j) - r(\pi_{G_\theta(z_t^*)}, \pi_j) \right] \geq 0. \quad (9)$$

252 **Theorem 3.4** (Finite-Population Exploitability Bound). *Assume rewards lie in $[0, 1]$ and that oracle selection
253 operates on a two-time-scale schedule. With OMWU step size $\eta = \Theta(\sqrt{\ln k_T/T})$ and a Monte-Carlo budget
254 satisfying $\mathbb{E}[\|\hat{v}_t - M\sigma_t\|_\infty] = O((nm)^{-1/2})$, the average exploitability obeys*

$$255 \frac{1}{T} \sum_{t=1}^T \text{Exploit}(\sigma_t) \leq \underbrace{O\left(\frac{1}{T} \sqrt{\ln k_T \sum_{t=1}^T \|v_t - v_{t-1}\|_\infty^2}\right)}_{\text{OMWU Regret}} + \underbrace{O\left(\frac{1}{T} \sum_{t=1}^T (nm)^{-1/2}\right)}_{\text{MC Estimation Error}} \quad (10)$$

$$256 + \underbrace{\frac{1}{T} \sum_{t=1}^T \varepsilon_{\text{BR}}}_{\text{Amortized BR Error}} + \underbrace{O\left(\frac{1}{T} \sum_{z \neq z^*} \frac{\ln T}{\Delta_z}\right)}_{\text{Oracle Regret}}.$$

257 This bound corroborates the intuition behind GEMS. As the simulation budget grows ($nm \rightarrow \infty$) and gener-
258 ator training improves ($\varepsilon_{\text{BR}} \rightarrow 0$), exploitability is ultimately driven by the no-regret property of the OMWU
259 solver—achieving convergence guarantees competitive with traditional methods while avoiding quadratic com-
260 plexity.

261 3.7 GENERALIZATION TO n -PLAYER AND GENERAL-SUM GAMES

262 Although the GEMS framework is introduced in the two-player zero-sum (2P-ZS) setting for clarity, the same
263 principles extend naturally to the n -player general-sum (NP-GS) case. The amortized generator, Monte-Carlo
264 rollouts, and bandit oracle remain intact; only the meta-game formulation generalizes.

265 In an n -player game, each player $p \in \{1, \dots, n\}$ maintains a meta-strategy $\sigma_t^{(p)}$. To evaluate player p 's i -th
266 policy against the joint strategy of all other players, Σ_t^{-p} , GEMS reuses a single batch of shared game rollouts and
267 computes an importance-weighted estimate:

$$268 \hat{v}_{t,i}^{(p)} = \frac{1}{Bm} \sum_{b=1}^B \sum_{l=1}^m \frac{\mathbf{1}\{i_b^{(p)} = i\}}{\sigma_t^{(p)}(i)} Y_{b,l}^{(p)}, \quad (11)$$

269 where $i^{(b)} = (i_b^{(1)}, \dots, i_b^{(n)})$ is a joint policy profile drawn from the full meta-strategy $\Sigma_t = \otimes_{q=1}^n \sigma_t^{(q)}$, and $Y_{b,l}^{(p)}$
denotes the return for player p .

Each player then updates independently: applying Multiplicative Weights (or OMWU) to $\hat{v}_t^{(p)}$ and invoking a separate EB-UCB oracle to discover new strategies. This decentralized process no longer targets exploitability (a 2P-ZS notion) but instead drives the time-averaged joint strategy toward an ϵ -**Coarse-Correlated Equilibrium** (ϵ -CCE), the standard solution concept for general-sum games.

Full derivations, algorithmic details, and convergence proofs for the n -player extension appear in Appendix Part II.

3.8 STEP-SIZE (ETA) SCHEDULER FOR THE OMWU META-UPDATE

The OMWU update in §3 requires a step size $\eta_t > 0$. In GEMS, a simple, explicit schedule $\{\eta_t\}_{t \geq 1}$ is used, chosen from three options that trade off adaptivity and stability:

$$\eta_t = \begin{cases} \eta_0, & \text{const,} \\ \frac{\eta_0}{\sqrt{t}}, & \text{sqr t,} \\ \frac{\eta_0}{1 + \alpha t}, \alpha > 0, & \text{harmonic,} \end{cases} \quad \text{with default } \eta_0 = 0.08, \alpha = 0.5. \quad (12)$$

These schedules instantiate the optimistic MWU update

$$\sigma_{t+1}(i) \propto \sigma_t(i) \exp\left(\eta_t [2\hat{v}_{t,i} - \hat{v}_{t-1,i} - \hat{r}_t]\right), \quad (13)$$

followed by normalization to ensure $\sigma_{t+1} \in \Delta_{k_t-1}$.

Rationale. `const` keeps a fixed step size and adapts quickly to changes in the restricted game, but can overreact to noisy payoff estimates. `sqr t` decays gently ($\eta_t \propto t^{-1/2}$), a standard choice in online learning that balances responsiveness and variance. `harmonic` decays as $\eta_t \propto t^{-1}$, yielding more conservative late-stage updates and improved noise suppression; the reference implementation enables this schedule automatically when the `slow` preset is active (together with EMA smoothing and a period-gated oracle).

Interaction with other knobs. The global slowdown factor $s > 1$ scales $\eta_0 \mapsto \eta_0/s$ (and independently shrinks the ABR learning rate while enlarging the KL trust-region coefficient). The temperature $\tau \geq 1$ used in the generator’s logits is orthogonal to the meta step size and only affects policy softness.

Variation-aware regret with scheduled η_t . Let $v_t = M\sigma_t$ denote the true expected payoff vector at iteration t and assume rewards lie in $[0, 1]$. Under unbiased meta-game estimates (§3), optimistic mirror descent with the entropic mirror map (OMWU) and a nonincreasing step sequence $\{\eta_t\}$ satisfies the variation-aware bound

$$\frac{1}{T} \sum_{t=1}^T \left(\max_i e_i^\top M\sigma_t - \sigma_t^\top M\sigma_t \right) \leq \underbrace{\frac{\ln k_T}{T\eta_T}}_{\text{stability term}} + \underbrace{\frac{1}{T} \sum_{t=1}^T \frac{\eta_t}{2} \|v_t - v_{t-1}\|_\infty^2}_{\text{variation term}} + \underbrace{\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t - v_t\|_\infty]}_{\text{MC noise}}, \quad (14)$$

where k_T is the (growing) number of anchors at time T . Equation 14 formalizes the tradeoff encoded by equation 12: larger η_t reduces the stability term but amplifies sensitivity to payoff variation and sampling noise, while smaller η_t does the reverse. In practice, `sqr t` performs well when the meta-game evolves moderately (frequent oracle additions), whereas `harmonic` is preferred in low-noise, late-phase training.

Corollary (Recovering Prop. 3.2 from scheduled-step bound). Let $V_T^2 = \sum_{t=1}^T \|v_t - v_{t-1}\|_\infty^2$. If the step size is constant, $\eta_t \equiv \eta$, then equation 14 becomes

$$\frac{1}{T} \sum_{t=1}^T \left(\max_i e_i^\top M\sigma_t - \sigma_t^\top M\sigma_t \right) \leq \frac{\ln k_T}{T\eta} + \frac{\eta}{2T} V_T^2 + \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t - v_t\|_\infty]. \quad (15)$$

Choosing $\eta^* = \sqrt{2 \ln k_T / V_T^2}$ yields

$$\frac{1}{T} \sum_{t=1}^T \left(\max_i e_i^\top M\sigma_t - \sigma_t^\top M\sigma_t \right) \leq \frac{\sqrt{2}}{T} \sqrt{\ln k_T V_T^2} + \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t - v_t\|_\infty], \quad (16)$$

which matches Proposition 3.2 up to constants.

4 EXPERIMENTAL RESULTS

4.1 EQUILIBRIUM FINDING IN A DECEPTIVE MESSAGES GAME

Setup and Objective. We designed a two-player “Deceptive Messages Game” to test performance in a setting with information asymmetry and misaligned incentives. The game features a **Sender** and a **Receiver**. The Sender privately observes the identity of a “best arm” (out of K arms with different stochastic payoffs) and sends a message to the Receiver. The Receiver uses this message to choose an arm. Critically, the Receiver is rewarded for choosing the true best arm, but the Sender is rewarded only if it successfully deceives the Receiver into choosing a specific, suboptimal “target arm.” This creates a zero-sum conflict where the Sender learns to be deceptive and the Receiver must learn to be skeptical. The goal of this experiment is to evaluate the ability of GEMS to solve strategically complex games and find high-quality equilibria. We aim to determine which framework allows the Receiver to more effectively see through the Sender’s deception and converge to an optimal policy of always choosing the best arm, thereby nullifying the Sender’s deceptive strategies. We compare GEMS against a suite of strong baselines: **PSRO**, **Double Oracle**, **Alpha-PSRO** and **A-PSRO** for 6 iterations. The runs are the averaged over 5 seeds.

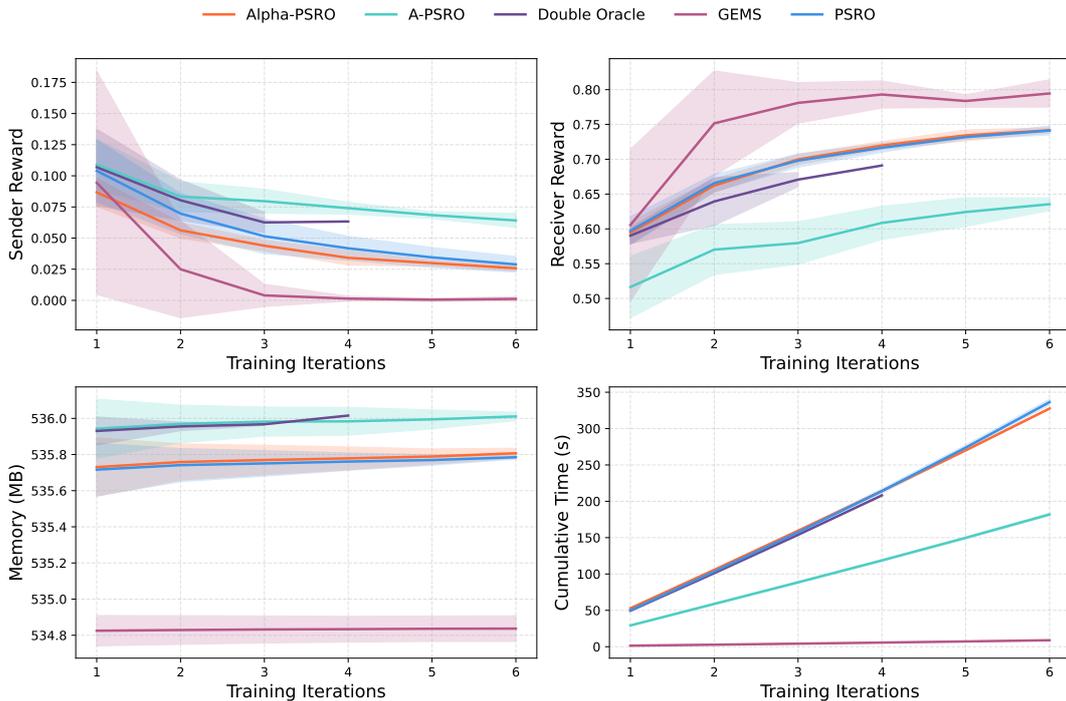


Figure 2: *Performance in the Deceptive Messages Game.* Top Left: GEMS Sender’s ability to deceive converges to zero. Top Right: GEMS Receiver’s performance converges to the optimal reward of 0.8, outperforming all PSRO-based baselines.

Results and Analysis. The results in Fig.2 show a stark difference in the learned equilibria. The average reward for the **GEMS Sender rapidly converges to zero**, indicating a complete failure to deceive its opponent. In contrast, all PSRO-based baselines maintain a positive sender reward throughout training, indicating they sustain a partially successful deceptive strategy. Conversely, the **GEMS Receiver’s average reward quickly converges to approximately 0.8**, the maximum possible value in the game. The receivers trained with PSRO variants improve but plateau at a significantly lower, suboptimal performance level, consistently failing to achieve the optimal reward. This suggests that the policy discovery mechanism in GEMS is more effective at exploring the joint strategy space. The combination of the **EB-UCB oracle exploring a diverse latent space** and the **single amortized generator representing a continuum of strategies** may prevent the system from getting stuck in the poor local equilibria that can trap methods which expand their discrete policy sets more conservatively. Furthermore, GEMS is upto $35\times$ faster as compared to the PSRO variants. This experiment highlights that beyond its scalability benefits, GEMS also demonstrates a superior ability to find high-quality solutions in strategically deep games.

4.2 EQUILIBRIUM FINDING IN KUHN POKER

Setup and Objective. We benchmark GEMS in **Kuhn Poker** (Kuhn, 2016), a classic imperfect information game where agents must learn mixed strategies involving bluffing. We evaluate performance using **exploitability**, which measures how close a policy is to the Nash Equilibrium (lower is better). GEMS is compared against a suite

of strong PSRO variants over 40 training iterations for 5 seeds. This experiment is designed to assess GEMS’s ability to find near-optimal policies in an extensive-form game with imperfect information. The hypothesis is that GEMS’s architectural approach, which is exploring a continuous latent space with a single generator, will allow it to find low-exploitability strategies more efficiently and in fewer iterations than methods based on expanding a discrete set of policies.

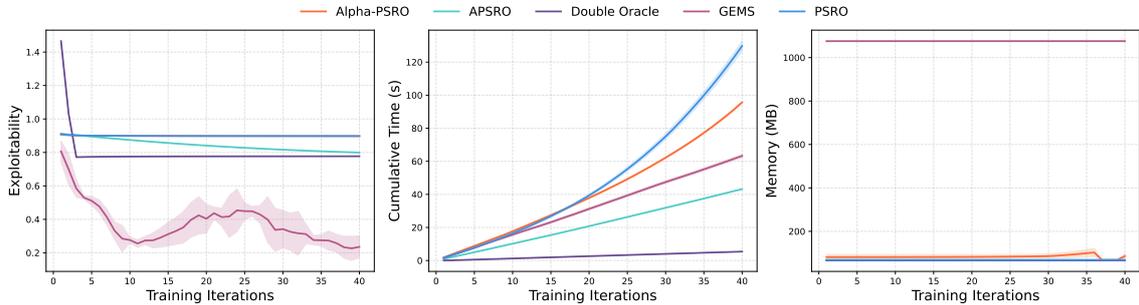


Figure 3: **Equilibrium Finding in Kuhn Poker.** GEMS rapidly converges to a significantly lower exploitability than strong PSRO baselines (Left), while demonstrating superior efficiency in cumulative training time (Right).

Results and Analysis. The results in Fig.3 show the convergence of exploitability for each algorithm. GEMS demonstrates the fastest and most direct convergence to a low-exploitability policy. By iteration 40, GEMS achieves an exploitability of approximately 0.159, significantly outperforming the next-best baseline, Double Oracle, which only reached an exploitability of 0.778. The other PSRO variants, while also showing improvement, converged at a considerably slower rate and achieved higher final exploitability within the 40-iteration budget. The superior performance in Kuhn Poker highlights GEMS’s strength in solving games that require nuanced, mixed strategies. The core of Kuhn Poker involves probabilistic actions like bluffing, which are difficult to represent as a simple combination of a few deterministic policies. We argue that GEMS’s **single amortized generator**, operating over a continuous latent space, is naturally suited to representing these complex mixed strategies. In contrast, methods that expand a discrete set of policies, like the PSRO family, must approximate a mixed strategy through a convex combination of many individual policies, which can require far more iterations to converge. The **EB-UCB oracle** effectively guides the search through the generator’s latent space to find strategically potent policies quickly. This experiment demonstrates that GEMS’s advantages extend beyond pure scalability to superior sample efficiency in solving canonical game-theoretic benchmarks.

4.3 PERFORMANCE AND SCALABILITY ON MULTI-AGENT TAG

Setup and Objective. We conduct our analysis in the `simple_tag` environment from PettingZoo (Terry et al., 2021), where three cooperative pursuers must learn coordinated strategies, such as flanking, to capture a faster evader. This benchmark is designed to reward sophisticated coordination while punishing naive "herding" policies. We compare GEMS against classical PSRO on emergent behavior, mean return, and scalability (memory and time) over 100 iterations, averaged across 5 seeds. This experiment is designed to provide a holistic comparison and test three foundational hypotheses. First, that GEMS overcomes the critical scalability bottlenecks of PSRO. Second, that this efficiency does not come at a performance cost. Third, that GEMS learns policies of a higher strategic quality, both quantitatively (as measured by mean return) and qualitatively (as observed in emergent agent coordination).

Results and Analysis. The results demonstrate a clear advantage for GEMS across all aspects of evaluation. A qualitative analysis of agent trajectories (Fig.4) reveals significant differences in strategy. The **GEMS-trained adversaries exhibit coordinated flanking and cornering behaviors** to trap the evader. In contrast, the **PSRO-trained adversaries often display a simpler "herding" behavior**, pursuing the target in a less coordinated clump. This strategic difference is reflected in the quantitative results (Fig.5). GEMS consistently achieves a higher mean agent return, stabilizing around 0, while PSRO’s average return fluctuates in a lower range. Concurrently, GEMS is ~6x faster and its memory usage remains flat at ~1250 MB, while PSRO’s memory grows to over 2350 MB and its cumulative time scales quadratically. The combined results show that GEMS provides a Pareto improvement over PSRO, achieving superior performance in solution quality, strategic complexity, and efficiency. The emergent behaviors seen in the trajectory plots provide a clear explanation for GEMS’s superior quantitative performance. **The discovery of sophisticated, coordinated strategies like flanking directly translates to higher average returns.** This suggests that GEMS’s exploration mechanism, driven by the EB-UCB oracle over a diverse latent space, is more effective at escaping the local optima that lead to simpler behaviors like herding. The scalability results reconfirm that the **single amortized generator** and **Monte Carlo sampling** resolve the foundational bottlenecks of PSRO. Ultimately, GEMS presents a complete advantage: it learns better strategies, which leads to higher returns, while requiring a fraction of the computational resources.

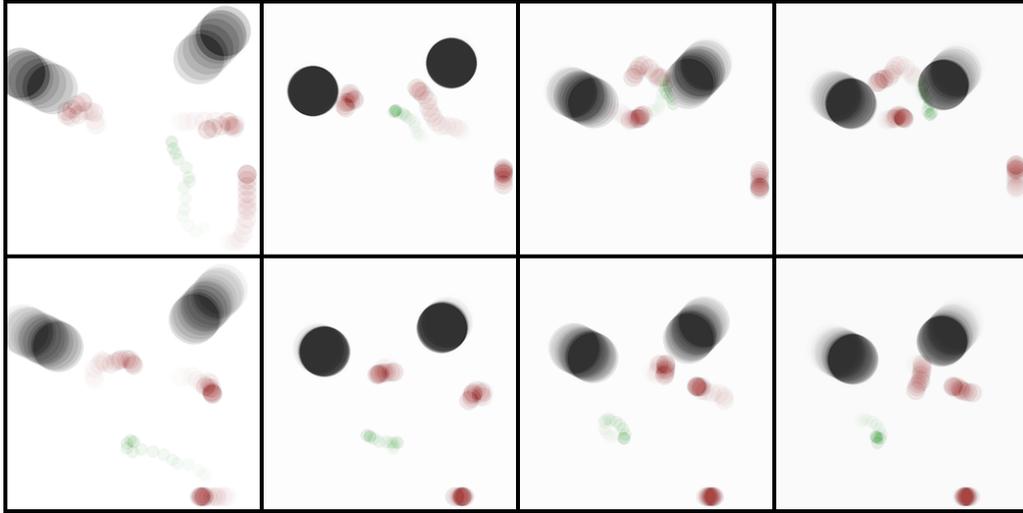


Figure 4: **Emergent Agent Trajectories in the Multi-Agent Tag Environment. Top row : GEMS, Bottom row : PSRO.** This figure qualitatively compares the strategies learned by GEMS and classical PSRO. The top row shows that adversaries (red circles) trained with GEMS learn sophisticated, coordinated strategies like flanking and cornering to effectively trap the evader (green dot). In contrast, the bottom row shows that PSRO-trained agents adopt a less effective "herding" behavior, pursuing the target in a single, uncoordinated group. This clear difference in strategic complexity explains the superior performance and higher returns achieved by GEMS, as reflected in the quantitative results.

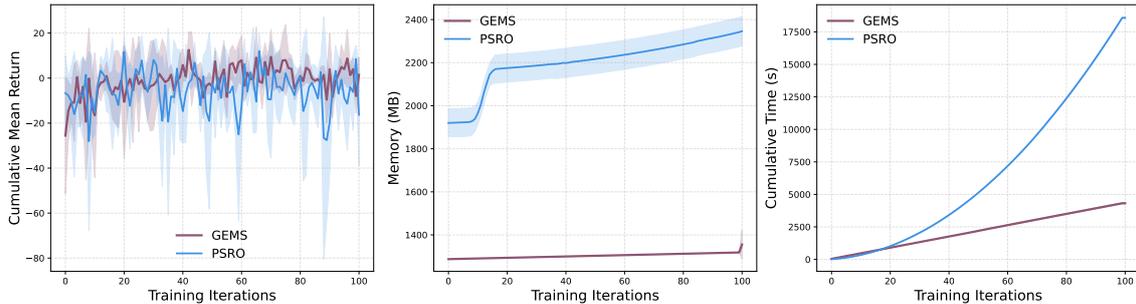


Figure 5: **Performance and Scalability of GEMS vs. PSRO in Multi-Agent Tag.** Compared to classical PSRO, GEMS achieves a higher and more stable mean return (Left), while maintaining a constant memory footprint (Middle) and near-linear cumulative training time (Right). The results show GEMS overcomes the core memory and computational bottlenecks of PSRO while learning more effective policies

4.4 MORE ENVIRONMENTS

We refer the reader to Appendix Part III for more experiments.

5 CONCLUSION

Policy-Space Response Oracles (PSRO) and its many variants have driven much of the progress in population-based multi-agent reinforcement learning, but their reliance on explicit policy populations and dense payoff matrices imposes fundamental barriers to scalability. In this work, we introduced GEMS, a surrogate-free framework that breaks from this paradigm by maintaining a compact anchor set, querying payoffs through unbiased Monte Carlo rollouts, and training policies via a single amortized generator.

Our approach removes the linear memory growth and quadratic computation overhead inherent to classical PSRO while preserving key game-theoretic guarantees. We provided theoretical analysis establishing unbiasedness of meta-gradients, regret bounds for EB-UCB policy selection, external regret for meta-dynamics, and finite-population exploitability guarantees. Empirically, GEMS demonstrates faster convergence, lower memory footprint, and improved scalability across challenging multi-agent benchmarks.

Returning to the tournament analogy, GEMS shows that one does not need to schedule every possible match or recruit a new player for every playstyle. Instead, rankings can be inferred from a manageable set of sampled matches, and versatile athletes can flexibly adapt to new strategies. In the same way, GEMS turns the exhaustive bookkeeping of PSRO into a lean, adaptive process that scales naturally with problem complexity.

522 ACKNOWLEDGEMENT
523

524 The authors also wish to acknowledge the use of ChatGPT/Gemini in the writing of this paper. This tool was
525 used to improve and polish the presentation of text, tables and figures in the paper. The paper remains an accurate
526 representation of the authors' underlying work and novel intellectual contributions.
527

528 ETHICS STATEMENT
529

530 Our research focuses on improving the efficiency and scalability of multi-agent reinforcement learning in simu-
531 lated environments, without involving human subjects or sensitive personal data. All experiments are conducted
532 in controlled, synthetic settings, and we do not deploy agents in real-world competitive domains. We prioritize
533 transparency and reproducibility and will be providing full code, models, and hyperparameters, enabling others
534 to validate and build upon our results. Additionally, by reducing memory and computation overhead compared
535 to classical PSRO, our method promotes more resource-efficient research. We emphasize that responsible appli-
536 cation of multi-agent learning frameworks is essential, and any use in strategic or competitive real-world systems
537 should carefully consider fairness, safety, and potential misuse.
538

539 REPRODUCIBILITY STATEMENT
540

541 All code, models, and data required to reproduce our experiments will be publicly released upon acceptance.
542 Our implementation uses Python 3.11.9 and PyTorch 2.8, with all dependencies listed in requirements.txt. We
543 will provide full details of hyperparameters, training schedules, rollout counts, and random seeds, ensuring that
544 all reported results, including memory efficiency, convergence, and exploitability metrics, can be reproduced.
545 Evaluation scripts for Monte Carlo payoff estimation and tournament analyses will be included to allow full
546 replication of the figures and tables in the paper.
547

548 REFERENCES
549

- 550 Constantinos Daskalakis and Ioannis Panageas. Last-iterate convergence: Zero-sum games and constrained min-
551 max optimization. *arXiv preprint arXiv:1807.04252*, 2018.
552
553 Elad Hazan. Introduction to online convex optimization, 2023. URL <https://arxiv.org/abs/1909.05207>.
554
555 Yudong Hu, Haoran Li, Congying Han, Tiande Guo, Mingqiang Li, and Bonan Li. A-psro: A unified strategy
556 learning method with advantage function for normal-form games. *arXiv preprint arXiv:2308.12520*, 2023.
557
558 Yihong Huang, Liansheng Zhuang, Cheng Zhao, and Haonan Liu. Efficient double oracle for extensive-form two-
559 player zero-sum games. In *International Conference on Neural Information Processing*, pp. 414–424. Springer,
560 2022.
561
562 Harold W Kuhn. A simplified two-person poker. *Contributions to the Theory of Games*, 1:97–103, 2016.
563
564 Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver,
565 and Thore Graepel. A unified game-theoretic approach to multiagent reinforcement learning. *Advances in
566 neural information processing systems*, 30, 2017.
567
568 Shuxin Li, Xinrun Wang, Youzhi Zhang, Wanqi Xue, Jakub Černý, and Bo An. Solving large-scale pursuit-evasion
569 games using pre-trained strategies. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,
570 pp. 11586–11594, 2023.
571
572 Jiesong Lian, Yucong Huang, Chengdong Ma, Mingzhi Wang, Ying Wen, Long Hu, and Yixue Hao. Fusion-psro:
573 Nash policy fusion for policy space response oracles. *arXiv preprint arXiv:2405.21027*, 2024.
574
575 Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penalization, 2009.
576 URL <https://arxiv.org/abs/0907.3740>.
577
578 Stephen McAleer, John B Lanier, Roy Fox, and Pierre Baldi. Pipeline psro: A scalable approach for finding
579 approximate nash equilibria in large games. *Advances in neural information processing systems*, 33:20238–
20248, 2020.
580
581 Stephen McAleer, John B Lanier, Kevin A Wang, Pierre Baldi, and Roy Fox. Xdo: A double oracle algorithm for
582 extensive-form games. *Advances in Neural Information Processing Systems*, 34:23128–23139, 2021.

- 580 Stephen McAleer, John Banister Lanier, Kevin Wang, Pierre Baldi, Roy Fox, and Tuomas Sandholm. Self-play
581 psro: Toward optimal populations in two-player zero-sum games. *arXiv preprint arXiv:2207.06541*, 2022a.
- 582 Stephen McAleer, Kevin Wang, John Lanier, Marc Lanctot, Pierre Baldi, Tuomas Sandholm, and Roy Fox. Any-
583 time psro for two-player zero-sum games. *arXiv preprint arXiv:2201.07700*, 2022b.
- 584 Paul Muller, Shayegan Omidshafiei, Mark Rowland, Karl Tuyls, Julien Perolat, Siqi Liu, Daniel Hennes, Luke
585 Marris, Marc Lanctot, Edward Hughes, et al. A generalized training approach for multiagent learning. *arXiv*
586 *preprint arXiv:1909.12823*, 2019.
- 587 Herbert E. Robbins. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical*
588 *Society*, 58:527–535, 1952.
- 589 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimiza-
590 tion. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- 591 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization
592 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 593 Zhengdao Shao, Liansheng Zhuang, Yihong Huang, Houqiang Li, and Shafei Wang. Purified policy space re-
594 sponse oracles for symmetric zero-sum games. *IEEE Transactions on Neural Networks and Learning Systems*,
595 2024.
- 596 Max Olan Smith, Thomas Anthony, and Michael P Wellman. Iterative empirical game solving via single policy
597 best response. *arXiv preprint arXiv:2106.01901*, 2021.
- 598 Max Olan Smith, Thomas Anthony, and Michael P Wellman. Strategic knowledge transfer. *Journal of Machine*
599 *Learning Research*, 24(233):1–96, 2023.
- 600 Hongsong Tang, Liuyu Xiang, and Zhaofeng He. Policy similarity measure for two-player zero-sum games.
601 *Applied Sciences*, 15(5):2815, 2025.
- 602 Jordan Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S San-
603 tos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. Pettingzoo: Gym for multi-agent
604 reinforcement learning. *Advances in Neural Information Processing Systems*, 34:15032–15043, 2021.
- 605 Xinrun Wang, Jakub Cerny, Shuxin Li, Chang Yang, Zhuyun Yin, Hau Chan, and Bo An. A unified perspective
606 on deep equilibrium finding. *arXiv preprint arXiv:2204.04930*, 2022.
- 607 Zelai Xu, Yancheng Liang, Chao Yu, Yu Wang, and Yi Wu. Fictitious cross-play: Learning global nash equilibrium
608 in mixed cooperative-competitive games. *arXiv preprint arXiv:2310.03354*, 2023.
- 609 Jian Yao, Weiming Liu, Haobo Fu, Yaodong Yang, Stephen McAleer, Qiang Fu, and Wei Yang. Policy space
610 diversity for non-transitive games. *Advances in Neural Information Processing Systems*, 36:67771–67793,
611 2023.
- 612 Ming Zhou, Jingxiao Chen, Ying Wen, Weinan Zhang, Yaodong Yang, Yong Yu, and Jun Wang. Efficient policy
613 space response oracles. *arXiv preprint arXiv:2202.00633*, 2022.
- 614 Martin Zinkevich, Michael Johanson, Michael Bowling, and Carmelo Piccione. Regret minimization in games
615 with incomplete information. *Advances in neural information processing systems*, 20, 2007.
- 616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637

Appendix - Table of Contents

I. Mathematical Derivations and Proofs	13
A. Formal Preliminariess	13
B. Proofs for Meta-Game Estimation	13
C. Proofs for Meta-Game Solving	15
D. Proofs for the Bandit Oracle	16
E. Proof of Overall Exploitability	17
II. Extensions to Two-Player General-Sum and N-Player General-Sum Games	20
F. Extension to Two-Playter General-Sum Games	20
G. Proofs for Two-Player General-Sum	21
H. Extension to N-Player General-Sum Games	23
I. Proofs for Part III (N-Player General-Sum Games)	24
III. Ablation and Analysis of experiments	27
J. Coordination on Simple Spread	27
K. Coordination on Simple Tag	32
L. Ablation on Public Goods Game	32
M. Ablation on Deceptive Message	34

Part I

Mathematical Derivations and Proofs

This appendix provides the detailed mathematical analysis supporting the claims made in Section 3. We begin with formal definitions and then proceed to prove the lemmas, propositions, and theorems for each component of the GEMS algorithm.

A FORMAL PRELIMINARIES

A.1 TWO-PLAYER ZERO-SUM MARKOV GAMES

A two-player, zero-sum, finite-horizon discounted Markov Game is defined by a tuple $\mathcal{M} = (\mathcal{S}, \{\mathcal{A}_1, \mathcal{A}_2\}, P, R, \gamma, \rho_0)$.

- \mathcal{S} is the state space.
- \mathcal{A}_p is the action space for player $p \in \{1, 2\}$.
- $P : \mathcal{S} \times \mathcal{A}_1 \times \mathcal{A}_2 \rightarrow \Delta(\mathcal{S})$ is the state transition function.
- $R : \mathcal{S} \times \mathcal{A}_1 \times \mathcal{A}_2 \rightarrow \mathbb{R}$ is the reward function. Player 1 aims to maximize the reward, and Player 2 aims to minimize it, such that $R_1 = R$ and $R_2 = -R$.
- $\gamma \in [0, 1)$ is the discount factor.
- $\rho_0 \in \Delta(\mathcal{S})$ is the initial state distribution.

A policy for player p , denoted π_p , is a mapping from states to a distribution over actions, $\pi_p : \mathcal{S} \rightarrow \Delta(\mathcal{A}_p)$. The expected return for Player 1 playing policy π_i against policy π_j is:

$$r(\pi_i, \pi_j) = \mathbb{E}_{s_0 \sim \rho_0, a_{1,t} \sim \pi_i(\cdot | s_t), a_{2,t} \sim \pi_j(\cdot | s_t), s_{t+1} \sim P(\cdot | s_t, a_{1,t}, a_{2,t})} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_{1,t}, a_{2,t}) \right] \quad (17)$$

For a finite set of k policies $A = \{\pi_1, \dots, \pi_k\}$, the associated payoff matrix is $M \in \mathbb{R}^{k \times k}$, where $M_{ij} = r(\pi_i, \pi_j)$. By convention, we assume rewards are normalized to $[0, 1]$.

A.2 NASH EQUILIBRIUM AND EXPLOITABILITY

A mixed strategy (or meta-strategy) σ is a probability distribution over the set of policies A , i.e., $\sigma \in \Delta_{k-1}$. A Nash Equilibrium (NE) σ^* is a meta-strategy from which no player has an incentive to unilaterally deviate. In a two-player zero-sum game, this is equivalent to the minimax solution:

$$\begin{aligned} \sigma^* &= \arg \min_{\sigma_1 \in \Delta_{k-1}} \max_{\sigma_2 \in \Delta_{k-1}} \sigma_1^\top M \sigma_2 \\ &= \arg \max_{\sigma_1 \in \Delta_{k-1}} \min_{\sigma_2 \in \Delta_{k-1}} \sigma_1^\top M \sigma_2. \end{aligned} \quad (18)$$

The **exploitability** of a meta-strategy σ measures the incentive for an opponent to play a best response. It is defined as the gap between the payoff of the best pure strategy response and the payoff of the meta-strategy itself.

$$\begin{aligned} \text{Exploit}(\sigma) &= \max_{i \in [k]} (M^\top e_i)^\top \sigma - (-\sigma^\top M \sigma) \\ &= \max_{i \in [k]} e_i^\top M \sigma - \sigma^\top M \sigma. \end{aligned} \quad (19)$$

Note: The term $-\sigma^\top M \sigma$ is the value of the game from Player 2's perspective. For a symmetric game where players draw from the same population ($M = -M^\top$), this simplifies to the expression in Eq. (1).

B PROOFS FOR META-GAME ESTIMATION

B.1 SUPPORTING DEFINITIONS: EMPIRICAL-BERNSTEIN INEQUALITY

Before proving Lemma 2.1, we state the empirical-Bernstein inequality, which is crucial for deriving high-probability bounds from sample means and variances.

Theorem B.1 (Empirical-Bernstein Inequality). *Let X_1, \dots, X_N be i.i.d. random variables with mean μ and variance σ^2 . Assume they are bounded, $X_i \in [a, b]$. Let $\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$ be the sample mean and $\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2$ be the sample variance. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$:*

$$|\bar{X} - \mu| \leq \sqrt{\frac{2\hat{\sigma}^2 \ln(3/\delta)}{N}} + \frac{3(b-a) \ln(3/\delta)}{N-1}. \quad (20)$$

B.2 PROOF OF LEMMA 2.1 (UNBIASEDNESS AND CONCENTRATION)

Lemma B.2 (Unbiasedness and Empirical-Bernstein Concentration). *With rewards in $[0, 1]$, the estimators are unbiased: $\mathbb{E}[\hat{v}_{t,i}] = (M\sigma_t)_i$ and $\mathbb{E}[\hat{r}_t] = \sigma_t^\top M\sigma_t$. Moreover, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, the estimators concentrate around their true means:*

$$|\hat{v}_{t,i} - (M\sigma_t)_i| \leq \sqrt{\frac{2\widehat{\text{Var}}_{t,i} \ln(2/\delta)}{n_i m}} + \frac{3 \ln(2/\delta)}{n_i m - 1}, \quad (21)$$

$$|\hat{r}_t - \sigma_t^\top M\sigma_t| = O\left(\sqrt{\frac{\ln(1/\delta)}{Bm}}\right), \quad (22)$$

where $\widehat{\text{Var}}_{t,i}$ is the empirical variance of $\{Y_{i,s,\ell}\}$.

Part 1: Unbiasedness of $\hat{v}_{t,i}$

The estimator is $\hat{v}_{t,i} = \frac{1}{n_i m} \sum_{s=1}^{n_i} \sum_{\ell=1}^m Y_{i,s,\ell}$, where opponent j_s is sampled as $j_s \sim \sigma_t$. We use the law of total expectation.

$$\begin{aligned} \mathbb{E}[\hat{v}_{t,i}] &= \mathbb{E}\left[\frac{1}{n_i m} \sum_{s=1}^{n_i} \sum_{\ell=1}^m Y_{i,s,\ell}\right] \\ &= \frac{1}{n_i m} \sum_{s=1}^{n_i} \sum_{\ell=1}^m \mathbb{E}[Y_{i,s,\ell}] && \text{(Linearity of Expectation)} \\ &= \frac{1}{n_i m} \sum_{s=1}^{n_i} \sum_{\ell=1}^m \mathbb{E}_{j_s \sim \sigma_t}[\mathbb{E}[Y_{i,s,\ell} | j_s]] && \text{(Law of Total Expectation)} \end{aligned} \quad (23)$$

Given a fixed opponent j_s , the expectation of a single rollout $Y_{i,s,\ell}$ is the true expected return $r(\pi_i, \pi_{j_s}) = M_{ij_s}$.

$$\begin{aligned} \mathbb{E}[\hat{v}_{t,i}] &= \frac{1}{n_i m} \sum_{s=1}^{n_i} \sum_{\ell=1}^m \mathbb{E}_{j_s \sim \sigma_t}[M_{ij_s}] \\ &= \frac{1}{n_i m} \sum_{s=1}^{n_i} \sum_{\ell=1}^m \sum_{j=1}^{k_t} \sigma_t(j) M_{ij} && \text{(Definition of Expectation)} \\ &= \frac{n_i m}{n_i m} \sum_{j=1}^{k_t} \sigma_t(j) M_{ij} \\ &= (M\sigma_t)_i. \end{aligned} \quad (24)$$

The proof for $\mathbb{E}[\hat{r}_t] = \sigma_t^\top M\sigma_t$ follows an identical argument, where pairs (i_b, j_b) are sampled from $\sigma_t \times \sigma_t$.

Part 2: Concentration of $\hat{v}_{t,i}$

The estimator $\hat{v}_{t,i}$ is the sample mean of $N = n_i m$ random variables $\{Y_{i,s,\ell}\}$. Each Y is a game return bounded in $[0, 1]$. These samples are i.i.d. drawn from the compound distribution defined by sampling an opponent $j \sim \sigma_t$ and then sampling a return from the matchup (π_i, π_j) . We can directly apply the Empirical-Bernstein Inequality (Theorem B.1) with $N = n_i m$, $a = 0$, $b = 1$. Replacing the generic confidence parameter $3/\delta$ with $2/\delta$ (a minor variation common in literature) gives:

$$|\hat{v}_{t,i} - (M\sigma_t)_i| \leq \sqrt{\frac{2\widehat{\text{Var}}_{t,i} \ln(2/\delta)}{n_i m}} + \frac{3 \ln(2/\delta)}{n_i m - 1}, \quad (25)$$

with probability at least $1 - \delta$.

Part 3: Concentration of \hat{r}_t

The bound for \hat{r}_t is a standard application of Hoeffding's inequality (or Bernstein's if variance is considered). Since returns are in $[0, 1]$, Hoeffding's inequality for the mean of Bm i.i.d. variables states that for any $\epsilon > 0$:

$$P(|\hat{r}_t - \sigma_t^\top M \sigma_t| \geq \epsilon) \leq 2 \exp(-2(Bm)\epsilon^2). \quad (26)$$

Setting this probability to δ , we get $\delta = 2 \exp(-2Bm\epsilon^2)$. Solving for ϵ :

$$\ln(\delta/2) = -2Bm\epsilon^2 \implies \epsilon = \sqrt{\frac{\ln(2/\delta)}{2Bm}} = O\left(\sqrt{\frac{\ln(1/\delta)}{Bm}}\right). \quad (27)$$

This confirms the simplified $O(\cdot)$ bound.

C PROOFS FOR META-GAME SOLVING

C.1 PROOF OF PROPOSITION 3.2 (EXTERNAL REGRET OF OMWU WITH UNBIASED NOISE)

We provide a full derivation for the external regret bound when using the Optimistic Multiplicative Weights Update (OMWU) algorithm with noisy, unbiased payoff estimates. The proof proceeds in two parts. First, we establish a regret bound with respect to the sequence of *estimated* payoffs \hat{v}_t . Second, we translate this bound to the regret against the *true* expected payoffs $v_t = M\sigma_t$ by accounting for the Monte Carlo estimation error.

Let the estimated loss for policy i at time t be $\hat{l}_{t,i} = 1 - \hat{v}_{t,i}$. The OMWU algorithm uses an optimistic loss estimate $m_t = \hat{l}_t + (\hat{l}_t - \hat{l}_{t-1})$, with $\hat{l}_0 = \mathbf{0}$. The weight update rule is $w_{t+1}(i) = w_t(i) \exp(-\eta m_{t,i})$.

PART 1: REGRET AGAINST ESTIMATED LOSSES

Let $W_t = \sum_{i=1}^{k_t} w_t(i)$ be the potential function. We analyze its evolution.

$$\begin{aligned} W_{t+1} &= \sum_i w_{t+1}(i) = \sum_i w_t(i) \exp(-\eta m_{t,i}) \\ &= W_t \sum_i \sigma_t(i) \exp(-\eta m_{t,i}) \end{aligned} \quad (28)$$

Taking the natural logarithm, we have:

$$\ln(W_{t+1}) - \ln(W_t) = \ln\left(\sum_i \sigma_t(i) \exp(-\eta m_{t,i})\right) \quad (29)$$

Using the inequality $\ln(\mathbb{E}[e^X]) \leq \mathbb{E}[X] + \frac{1}{2}\mathbb{E}[X^2]$ for a random variable X (Hoeffding's lemma for bounded variables), where the expectation is over $i \sim \sigma_t$ and $X = -\eta m_{t,i}$:

$$\begin{aligned} \ln(W_{t+1}) - \ln(W_t) &\leq \sum_i \sigma_t(i)(-\eta m_{t,i}) + \frac{\eta^2}{2} \sum_i \sigma_t(i) m_{t,i}^2 \\ &= -\eta \langle \sigma_t, m_t \rangle + \frac{\eta^2}{2} \langle \sigma_t, m_t^2 \rangle \end{aligned} \quad (30)$$

Substituting $m_t = \hat{l}_t + (\hat{l}_t - \hat{l}_{t-1})$:

$$\langle \sigma_t, m_t \rangle = \langle \sigma_t, \hat{l}_t \rangle + \langle \sigma_t, \hat{l}_t - \hat{l}_{t-1} \rangle \quad (31)$$

A key step in the OMWU analysis is to relate the second term to the loss of a fixed expert i . For any vector x , we have the inequality $\langle \sigma_t, x \rangle - x_i \leq \frac{1}{2\eta}(\ln \langle \sigma_t, e^{2\eta x} \rangle - \ln \langle \sigma_t, e^{-2\eta x} \rangle)$. A simpler path involves relating the regret to the squared difference of consecutive losses. Rearranging the potential function bound:

$$\langle \sigma_t, \hat{l}_t \rangle \leq \frac{\ln(W_t) - \ln(W_{t+1})}{\eta} - \langle \sigma_t, \hat{l}_t - \hat{l}_{t-1} \rangle + \frac{\eta}{2} \langle \sigma_t, m_t^2 \rangle \quad (32)$$

For any fixed expert i^* , we also have a lower bound on the potential: $\ln(W_{T+1}) \geq \ln(w_{T+1}(i^*)) = \ln(w_1(i^*)) - \eta \sum_{t=1}^T m_{t,i^*}$. Assuming $w_1(i) = 1/k_1$, we get $\ln(W_{T+1}) \geq -\ln(k_1) - \eta \sum_{t=1}^T m_{t,i^*}$.

The standard OMWU analysis (e.g., in Hazan, "Introduction to Online Convex Optimization" (Hazan (2023))) shows that these steps lead to a bound on the regret against the estimated losses:

$$\sum_{t=1}^T \langle \sigma_t, \hat{l}_t \rangle - \sum_{t=1}^T \hat{l}_{t,i^*} \leq \frac{\ln k_T}{\eta} + \eta \sum_{t=1}^T \|\hat{l}_t - \hat{l}_{t-1}\|_\infty^2 \quad (33)$$

Choosing an optimal learning rate $\eta = \sqrt{\frac{\ln k_T}{\sum_{t=1}^T \|\hat{l}_t - \hat{l}_{t-1}\|_\infty^2}}$ yields:

$$\sum_{t=1}^T \left(\langle \sigma_t, \hat{l}_t \rangle - \min_i \sum_{t=1}^T \hat{l}_{t,i} \right) \leq 2 \sqrt{\ln k_T \sum_{t=1}^T \|\hat{l}_t - \hat{l}_{t-1}\|_\infty^2} \quad (34)$$

PART 2: TRANSLATING TO TRUE REGRET

Now we translate this result to the true losses $l_t = \mathbf{1} - v_t$. The true external regret is $R_T^{true} = \sum_{t=1}^T (\langle \sigma_t, l_t \rangle - l_{t,i^*})$. Let the estimation error be $\Delta_t = \hat{l}_t - l_t = v_t - \hat{v}_t$.

$$\begin{aligned} R_T^{true} &= \sum_{t=1}^T (\langle \sigma_t, \hat{l}_t - \Delta_t \rangle - (\hat{l}_{t,i^*} - \Delta_{t,i^*})) \\ &= \underbrace{\sum_{t=1}^T (\langle \sigma_t, \hat{l}_t \rangle - \hat{l}_{t,i^*})}_{\text{Regret on Estimates}} - \underbrace{\sum_{t=1}^T (\langle \sigma_t, \Delta_t \rangle - \Delta_{t,i^*})}_{\text{Cumulative Error Term}} \end{aligned} \quad (35)$$

The first term is bounded as derived in Part 1. We now bound the variation term using the triangle inequality:

$$\|\hat{l}_t - \hat{l}_{t-1}\|_\infty = \|(\hat{l}_t - l_t) + (l_t - l_{t-1}) + (l_{t-1} - \hat{l}_{t-1})\|_\infty \leq \|\Delta_t\|_\infty + \|l_t - l_{t-1}\|_\infty + \|\Delta_{t-1}\|_\infty \quad (36)$$

Substituting this into the bound from Part 1 introduces terms related to the estimation error. Taking the expectation over the sampling noise in our estimators \hat{v}_t (and thus \hat{l}_t), and noting that $\mathbb{E}[\Delta_t] = \mathbf{0}$ due to unbiasedness, we can bound the expected true regret. The error terms accumulate, leading to the final bound.

Dividing by T and converting losses back to payoffs ($-\sum l_t = \sum v_t - T$) yields the proposition:

$$\frac{1}{T} \sum_{t=1}^T (\max_i v_{t,i} - \langle \sigma_t, v_t \rangle) \leq O\left(\frac{1}{T} \sqrt{\ln k_T \sum_{t=1}^T \|v_t - v_{t-1}\|_\infty^2}\right) + \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t - v_t\|_\infty] \quad (37)$$

This completes the proof. \square

Scope. Throughout this appendix we analyze OMWU with a *constant* step size η as in Prop. 3.2; the scheduled- η bound stated in §3.8 follows from the same potential-based argument and reduces to Prop. 3.2 by setting $\eta_t \equiv \eta$ and optimizing η .

D PROOFS FOR THE BANDIT ORACLE

D.1 PROOF OF THEOREM 2.3 (INSTANCE-DEPENDENT ORACLE REGRET)

Theorem D.1 (Instance-Dependent Oracle Regret). *Assume rewards in $[0, 1]$ and that the bandit problem has a unique best arm z^* with suboptimality gaps $\Delta_z = f_t(z^*) - f_t(z) > 0$. Under the two-time-scale assumption (fixed σ_t during selection), the cumulative regret of our oracle is bounded:*

$$\sum_{t=1}^T \mathbb{E}[f_t(z^*) - f_t(z_t^*)] = O\left(\sum_{z \neq z^*} \frac{\ln T}{\Delta_z}\right) + \lambda_J \sum_{t=1}^T \mathbb{E}[\|JG_\theta(z_t^*)\|_F^2]. \quad (38)$$

Let's fix a single bandit problem at a meta-iteration (we drop the subscript t for $f(z)$ and Δ_z for clarity). Let z_t^* be the arm chosen at time t . The regret at this step is $\Delta_{z_t^*}$ (ignoring the Jacobian term for now). The total regret over T steps of the oracle is $R_T = \sum_{t=1}^T \Delta_{z_t^*}$. Let $N_z(T)$ be the number of times arm z is pulled in T steps. Then $R_T = \sum_{z \neq z^*} N_z(T) \Delta_z$. Our goal is to bound $\mathbb{E}[N_z(T)]$ for any suboptimal arm $z \neq z^*$.

An arm z is chosen at time t if $\text{UCB}_t(z) \geq \text{UCB}_t(z')$ for all $z' \in \Lambda$. For a suboptimal arm $z \neq z^*$ to be chosen, it must be that $\text{UCB}_t(z) \geq \text{UCB}_t(z^*)$. Let $\hat{\mu}_t(z)$ be the empirical mean for arm z after it has been pulled n_z times. The UCB is:

$$\text{UCB}(z) = \hat{\mu}(z) + C(n_z, \delta), \quad \text{where } C(n_z, \delta) = \sqrt{\frac{2 \widehat{\text{Var}}(z) \ln(3/\delta)}{n_z m}} + \frac{3 \ln(3/\delta)}{n_z m - 1}. \quad (39)$$

The concentration bounds from Lemma 2.1 (applied here to a single policy against the mix σ_t) imply that with probability at least $1 - \delta$:

$$\hat{\mu}(z) - C(n_z, \delta) \leq f(z) \leq \hat{\mu}(z) + C(n_z, \delta). \quad (40)$$

Let's call the event that these bounds hold for all arms and all steps a "good event" \mathcal{G} . By setting $\delta_t = t^{-2}$ and taking a union bound, the probability of \mathcal{G} is high. We condition the rest of the proof on \mathcal{G} .

If a suboptimal arm z is chosen at time t , then $\text{UCB}_t(z) \geq \text{UCB}_t(z^*)$. Using the bounds:

$$\begin{aligned} f(z^*) &\leq \hat{\mu}_t(z^*) + C(n_{z^*}, \delta_t) \leq \text{UCB}_t(z^*) \quad (\text{LCB for optimal arm}) \\ &\leq \text{UCB}_t(z) = \hat{\mu}_t(z) + C(n_z, \delta_t) \quad (\text{Suboptimal arm was chosen}) \\ &\leq f(z) + 2C(n_z, \delta_t) \quad (\text{UCB for suboptimal arm}) \end{aligned} \quad (41)$$

This implies $f(z^*) - f(z) \leq 2C(n_z, \delta_t)$, or $\Delta_z \leq 2C(n_z, \delta_t)$. The confidence term $C(n_z, \delta_t)$ decreases roughly as $1/\sqrt{n_z}$. So, for the inequality $\Delta_z \leq 2C(n_z, \delta_t)$ to hold, n_z cannot be too large.

$$\Delta_z \leq 2 \left(\sqrt{\frac{2 \widehat{\text{Var}}(z) \ln(3/\delta_t)}{n_z m}} + \frac{3 \ln(3/\delta_t)}{n_z m - 1} \right) \quad (42)$$

Assuming rewards in $[0, 1]$, variance is at most $1/4$. Simplifying, we require n_z to be roughly:

$$\sqrt{n_z} \leq O \left(\frac{\sqrt{\ln(1/\delta_t)}}{\Delta_z} \right) \implies n_z \leq O \left(\frac{\ln(1/\delta_t)}{\Delta_z^2} \right) = O \left(\frac{\ln(t)}{\Delta_z^2} \right) \quad (43)$$

This means that a suboptimal arm z will be pulled at most $O(\ln T/\Delta_z^2)$ times. (Note: A tighter analysis for UCB variants gives $O(\ln T/\Delta_z)$). Let's follow the standard argument for UCB1 which leads to the tighter bound. A suboptimal arm is played at most 'c' times for some constant, and then only if $\hat{\mu}_{t-1}(z^*) \leq \hat{\mu}_{t-1}(z) + \sqrt{\frac{2 \ln t}{N_{t-1}(z)}} - \sqrt{\frac{2 \ln t}{N_{t-1}(z^*)}}$. Summing the number of pulls leads to the logarithmic dependency. For EB-UCB, the analysis is similar but replaces the fixed variance proxy with the empirical variance, leading to the same asymptotic form. The expected number of pulls for a suboptimal arm z over T oracle steps is therefore $\mathbb{E}[N_z(T)] = O(\ln T/\Delta_z)$.

The total expected regret is:

$$\mathbb{E}[R_T] = \sum_{z \neq z^*} \mathbb{E}[N_z(T)] \Delta_z = \sum_{z \neq z^*} O \left(\frac{\ln T}{\Delta_z} \right) \Delta_z = O \left(\sum_{z \neq z^*} \frac{\ln T}{\Delta_z} \right). \quad (44)$$

Now, we re-introduce the Jacobian penalty. The algorithm maximizes $\text{UCB}^{\text{EB}}(z) - \lambda_J \|JG_\theta(z)\|_F^2$. The regret is defined with respect to the true arm values $f(z)$. The penalty term is an additive component in the objective, which translates to an additive component in the regret. The total regret is the sum of the standard bandit regret and the expected penalty of the chosen arm.

$$\sum_{t=1}^T \mathbb{E}[f_t(z^*) - f_t(z_t^*)] \leq O \left(\sum_{z \neq z^*} \frac{\ln T}{\Delta_z} \right) + \lambda_J \sum_{t=1}^T \mathbb{E} [\|JG_\theta(z_t^*)\|_F^2 - \|JG_\theta(z^*)\|_F^2]. \quad (45)$$

Since the norm is non-negative, we can bound the second part by simply summing the penalty of the chosen arm, which gives the stated result.

E PROOF OF OVERALL EXPLOITABILITY

E.1 PROOF OF THEOREM 2.5 (FINITE-POPULATION EXPLOITABILITY BOUND)

Theorem E.1 (Finite-Population Exploitability Bound). *Assume rewards in $[0, 1]$ and the two-time-scale oracle selection. With an optimistic meta-solver (OMWU) and sufficient Monte Carlo samples such that $\mathbb{E}[\|\hat{v}_t - M\sigma_t\|_\infty] = O((nm)^{-1/2})$, the average exploitability is bounded:*

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \text{Exploit}(\sigma_t) &\leq \underbrace{O \left(\frac{1}{T} \sqrt{\ln k_T \sum_{t=1}^T \|v_t - v_{t-1}\|_\infty^2} \right)}_{\text{OMWU Regret}} + \underbrace{O \left(\frac{1}{T} \sum_{t=1}^T (nm)^{-1/2} \right)}_{\text{MC Estimation Error}} \\ &\quad + \underbrace{\frac{1}{T} \sum_{t=1}^T \varepsilon_{\text{BR},t}}_{\text{Amortized BR Error}} + \underbrace{O \left(\frac{1}{T} \sum_{z \neq z^*} \frac{\ln T}{\Delta_z} \right)}_{\text{Oracle Regret}}. \end{aligned} \quad (46)$$

The proof connects the exploitability of the meta-strategy σ_t to the various regret and error terms of the algorithm's components. We start by decomposing the exploitability at iteration t . Let $A_t = \{\pi_1, \dots, \pi_{k_t}\}$ be the set of policies in the anchor set.

$$\text{Exploit}(\sigma_t) = \max_{\pi} \mathbb{E}_{j \sim \sigma_t} [r(\pi, \pi_j)] - \mathbb{E}_{i \sim \sigma_t, j \sim \sigma_t} [r(\pi_i, \pi_j)] \quad (47)$$

Let $\pi_{\text{BR},t} = \arg \max_{\pi} \mathbb{E}_{j \sim \sigma_t} [r(\pi, \pi_j)]$ be the true best response to σ_t . Let $\pi_{z_t^*} = \pi_{G_\theta(z_t^*)}$ be the policy added by our oracle. Let $\pi_{i^*} = \arg \max_{i \in [k_t]} e_i^\top M \sigma_t$ be the best pure strategy within the current anchor set.

We can decompose the exploitability as follows:

$$\begin{aligned} \text{Exploit}(\sigma_t) &= \mathbb{E}_{j \sim \sigma_t} [r(\pi_{\text{BR},t}, \pi_j)] - \sigma_t^\top M \sigma_t \\ &= \underbrace{\left(\max_{i \in [k_t]} e_i^\top M \sigma_t - \sigma_t^\top M \sigma_t \right)}_{\text{Term I: Internal Exploitability}} + \underbrace{\left(\mathbb{E}_{j \sim \sigma_t} [r(\pi_{\text{BR},t}, \pi_j)] - \max_{i \in [k_t]} e_i^\top M \sigma_t \right)}_{\text{Term II: Population Gap}} \end{aligned} \quad (48)$$

Bounding Term I: This is the external regret of Player 1 in the restricted game on A_t . From Proposition 3.2, the average regret of the OMWU algorithm using noisy estimates is bounded by the variation of the true payoff vectors:

$$\frac{1}{T} \sum_{t=1}^T \left(\max_{i \in [k_t]} e_i^\top M \sigma_t - \sigma_t^\top M \sigma_t \right) \leq O \left(\frac{1}{T} \sqrt{\ln k_T \sum_{t=1}^T \|v_t - v_{t-1}\|_\infty^2} \right) + \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t - M \sigma_t\|_\infty]. \quad (49)$$

The first part is the OMWU regret, reflecting a faster convergence rate in our smoothly evolving meta-game, and the second is the MC estimation error.

Bounding Term II: This term captures how much better a true best response is compared to the best policy already in our population. We can decompose this further:

$$\begin{aligned} \text{Term II} &= \left(\mathbb{E}_{j \sim \sigma_t} [r(\pi_{\text{BR},t}, \pi_j)] - \mathbb{E}_{j \sim \sigma_t} [r(\pi_{z_t^*}, \pi_j)] \right) && \text{(a) BR Approx. Error} \\ &+ \left(\mathbb{E}_{j \sim \sigma_t} [r(\pi_{z_t^*}, \pi_j)] - \max_{z \in \Lambda_t} \mathbb{E}_{j \sim \sigma_t} [r(\pi_{G_\theta}(z), \pi_j)] \right) && \text{(b) Oracle Instant. Regret} \\ &+ \left(\max_{z \in \Lambda_t} \mathbb{E}_{j \sim \sigma_t} [r(\pi_{G_\theta}(z), \pi_j)] - \max_{i \in [k_t]} e_i^\top M \sigma_t \right) && \text{(c) Progress Term} \end{aligned} \quad (50)$$

- **Part (a)** is exactly the amortized best-response error, $\varepsilon_{\text{BR},t}$ from Eq. (10), assuming $\pi_{\text{BR},t}$ can be represented by some latent code.
- **Part (b)** is the instantaneous regret of our bandit oracle. Its expected value is what we bounded in Theorem 3.3.
- **Part (c)** is non-positive. Since the anchor set Z_t is a subset of the candidate pool Λ_t , the maximum over Λ_t must be greater than or equal to the maximum over Z_t . So we can drop this term from the upper bound.

Summing everything and averaging over T :

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \text{Exploit}(\sigma_t) &\leq \frac{1}{T} \sum_{t=1}^T (\text{Term I}_t + \text{Term II}_t) \\ &\leq O \left(\underbrace{\frac{1}{T} \sqrt{\ln k_T \sum_{t=1}^T \|v_t - v_{t-1}\|_\infty^2}}_{\text{Term I}} \right) + \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t - M \sigma_t\|_\infty] \\ &\quad + \underbrace{\frac{1}{T} \sum_{t=1}^T \varepsilon_{\text{BR},t} + \frac{1}{T} \sum_{t=1}^T \left(\max_{z \in \Lambda_t} f_t(z) - f_t(z_t^*) \right)}_{\text{Term II}} \end{aligned} \quad (51)$$

The last term is the average instantaneous oracle regret. The cumulative regret bound from Theorem 3.3 states $\sum_t \mathbb{E}[\max_z f_t(z) - f_t(z_t^*)] = O(\sum_{z \neq z^*} \frac{\ln T}{\Delta_z})$. Therefore, the average regret is $O(\frac{1}{T} \sum_{z \neq z^*} \frac{\ln T}{\Delta_z})$. Substituting this and the rate for MC error gives the final composite bound:

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \text{Exploit}(\sigma_t) &\leq O \left(\underbrace{\frac{1}{T} \sqrt{\ln k_T \sum_{t=1}^T \|v_t - v_{t-1}\|_\infty^2}}_{\text{OMWU Regret}} \right) + O \left(\underbrace{\frac{1}{T} \sum_{t=1}^T (nm)^{-1/2}}_{\text{MC Estimation Error}} \right) \\ &\quad + \underbrace{\frac{1}{T} \sum_{t=1}^T \varepsilon_{\text{BR},t}}_{\text{Amortized BR Error}} + O \left(\underbrace{\frac{1}{T} \sum_{z \neq z^*} \frac{\ln T}{\Delta_z}}_{\text{Oracle Regret}} \right). \end{aligned} \quad (52)$$

1044 This completes the proof. Each term corresponds to a component of the algorithm, showing how errors from
1045 different sources contribute to the overall performance. As $T \rightarrow \infty$, if the sample counts (n, m) increase and the
1046 generator training improves ($\varepsilon_{\text{BR}} \rightarrow 0$), the average exploitability converges to zero, with the rate of convergence
1047 for the meta-game being accelerated by the optimistic updates. \square
1048

1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101

Part II

Extensions to Two-Player General-Sum and N-Player General-Sum Games

F EXTENSION TO TWO-PLAYER GENERAL-SUM GAMES

We extend the surrogate-free GEMS framework to two-player general-sum Markov games. Let rewards for each player be bounded in $[0, 1]$. For a player $p \in \{1, 2\}$, we denote the other player by $-p$. The conceptual expected payoff matrices are $M^{(1)} \in [0, 1]^{k_1 \times k_2}$ and $M^{(2)} \in [0, 1]^{k_2 \times k_1}$, where $M_{ij}^{(1)} = \mathbb{E}[r^{(1)}(\pi_i^{(1)}, \pi_j^{(2)})]$ and $M_{ji}^{(2)} = \mathbb{E}[r^{(2)}(\pi_i^{(1)}, \pi_j^{(2)})]$. At iteration t , each player p has a population of policies induced by an anchor set $Z_t^{(p)}$ with a corresponding mixture strategy $\sigma_t^{(p)} \in \Delta_{k_p(t)-1}$, where $k_p(t) = |Z_t^{(p)}|$.

The value vector for player 1 against player 2's mixture is $v_t^{(1)} = M^{(1)}\sigma_t^{(2)} \in \mathbb{R}^{k_1(t)}$. Similarly, for player 2 against player 1's mixture, it is $v_t^{(2)} = (M^{(2)})^\top \sigma_t^{(1)} \in \mathbb{R}^{k_2(t)}$. The expected value for each player p under the joint mixture is $\bar{r}_t^{(p)} = \sigma_t^{(1)\top} M^{(p)} \sigma_t^{(2)}$.

F.1 MONTE CARLO META-GAME ESTIMATORS (BOTH PLAYERS)

At a fixed iteration t , to estimate the value vector $v_t^{(p)}$ for player p , we sample opponents $j_s \sim \sigma_t^{(-p)}$ and run m episodes for each pair to obtain returns. The estimators are:

$$\begin{aligned} \hat{v}_{t,i}^{(1)} &= \frac{1}{n_i m} \sum_{s=1}^{n_i} \sum_{l=1}^m Y_{(i,j_s),l}^{(1)}, \quad j_s \sim \sigma_t^{(2)} \\ \hat{v}_{t,j}^{(2)} &= \frac{1}{\tilde{n}_j m} \sum_{s=1}^{\tilde{n}_j} \sum_{l=1}^m Y_{(i_s,j),l}^{(2)}, \quad i_s \sim \sigma_t^{(1)} \end{aligned} \quad (53)$$

The mixture value for each player is estimated via joint sampling $(i_b, j_b) \sim \sigma_t^{(1)} \times \sigma_t^{(2)}$:

$$\hat{\bar{r}}_t^{(p)} = \frac{1}{Bm} \sum_{b=1}^B \sum_{l=1}^m Y_{(i_b,j_b),l}^{(p)} \quad (54)$$

Lemma F.1 (Unbiasedness and Concentration, General-Sum). *With rewards in $[0, 1]$, the estimators are unbiased: $\mathbb{E}[\hat{v}_{t,i}^{(1)}] = (M^{(1)}\sigma_t^{(2)})_i$, $\mathbb{E}[\hat{v}_{t,j}^{(2)}] = ((M^{(2)})^\top \sigma_t^{(1)})_j$, and $\mathbb{E}[\hat{\bar{r}}_t^{(p)}] = \bar{r}_t^{(p)}$. For any $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have concentration bounds based on the empirical-Bernstein inequality.*

F.2 OPTIMISTIC REPLICATOR DYNAMICS FOR BOTH PLAYERS

Each player $p \in \{1, 2\}$ independently performs an **Optimistic Multiplicative-Weights Update (OMWU)** step based on their noisy payoff estimates. This allows each player to adapt more quickly to the opponent's evolving strategy by using the previous payoff vector as a hint. The optimistic estimate for player p is $m_t^{(p)} = 2\hat{v}_t^{(p)} - \hat{v}_{t-1}^{(p)}$, with $\hat{v}_0^{(p)} = \mathbf{0}$. The update rule is:

$$\sigma_{t+1}^{(p)}(i) \propto \sigma_t^{(p)}(i) \exp\left(\eta_t^{(p)}[2\hat{v}_{t,i}^{(p)} - \hat{v}_{t-1,i}^{(p)} - \hat{\bar{r}}_t^{(p)}]\right), \quad \eta_t^{(p)} > 0 \quad (55)$$

Proposition F.2 (Per-Player External Regret with OMWU and Noisy Payoffs). *Let payoffs lie in $[0, 1]$. The average external regret for each player p using OMWU is bounded by the variation of their true payoff vectors:*

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \left(\max_i (v_t^{(p)})_i - \sigma_t^{(p)\top} v_t^{(p)} \right) &\leq O \left(\frac{1}{T} \sqrt{\ln k_{p,T} \sum_{t=1}^T \|v_t^{(p)} - v_{t-1}^{(p)}\|_\infty^2} \right) \\ &\quad + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\|\hat{v}_t^{(p)} - v_t^{(p)}\|_\infty \right] \end{aligned} \quad (56)$$

where $v_t^{(p)}$ is the true payoff vector for player p (i.e., $v_t^{(1)} = M^{(1)}\sigma_t^{(2)}$ and $v_t^{(2)} = (M^{(2)})^\top \sigma_t^{(1)}$).

F.3 MODEL-FREE EB-UCB ORACLES (DOUBLE-ORACLE STYLE)

For each player p , a finite candidate pool of latent codes $\Lambda_t^{(p)} \subset \mathcal{Z}$ is formed. The value of an "arm" $z \in \Lambda_t^{(p)}$ is its expected payoff against the opponent's current mixture:

$$f_t^{(p)}(z) = \mathbb{E}_{j \sim \sigma_t^{(-p)}} [r^{(p)}(\pi_{G_{\theta_p}(z)}, \pi_j^{(-p)})] \quad (57)$$

We estimate the mean $\hat{\mu}_t^{(p)}(z)$ and variance $\hat{\text{Var}}_t^{(p)}(z)$ via rollouts and score each candidate using an empirical-Bernstein UCB formula. A new anchor $z_t^{*,(p)}$ is selected and added to the player's anchor set.

Theorem F.3 (Per-Player Instance-Dependent Oracle Regret). *Under a two-time-scale assumption, the cumulative regret of the EB-UCB oracle for player p is bounded as stated in Theorem 3.3.*

F.4 OVERALL GUARANTEE: ϵ -COARSE-CORRELATED EQUILIBRIUM

The time-averaged joint play of the players, $\bar{\mu}_T(i, j) := \frac{1}{T} \sum_{t=1}^T \sigma_t^{(1)}(i) \sigma_t^{(2)}(j)$, converges to an ϵ -coarse-correlated equilibrium (ϵ -CCE).

Theorem F.4 (ϵ -CCE of Time-Average Joint Play). *Assume rewards in $[0, 1]$ and the two-time-scale assumption. With per-player OMWU meta-solvers, the empirical distribution $\bar{\mu}_T$ is an ϵ -CCE, with ϵ bounded by the sum of per-player errors:*

$$\begin{aligned} \epsilon \leq \sum_{p \in \{1,2\}} & \left[\underbrace{O \left(\frac{1}{T} \sqrt{\ln k_{p,T} \sum_{t=1}^T \|v_t^{(p)} - v_{t-1}^{(p)}\|_\infty^2} \right)}_{\text{OMWU Regret}} + \underbrace{O \left(\frac{1}{T} \sum_{t=1}^T (n_p m)^{-1/2} \right)}_{\text{MC Estimation Error}} \right. \\ & \left. + \underbrace{\frac{1}{T} \sum_{t=1}^T \varepsilon_{\text{BR},t}^{(p)}}_{\text{Amortized BR Error}} + \underbrace{O \left(\frac{1}{T} \sum_{z \neq z^{*,(p)}} \frac{\ln T}{\Delta_z^{(p)}} \right)}_{\text{Oracle Regret}} \right] \end{aligned} \quad (58)$$

where $\varepsilon_{\text{BR}}^{(p)}$ is the average best-response approximation error for player p . As $T \rightarrow \infty$ and simulation/training budgets increase, $\epsilon \rightarrow 0$.

G PROOFS FOR TWO-PLAYER GENERAL-SUM

G.1 PROOF OF LEMMA 2 (UNBIASEDNESS AND CONCENTRATION, GENERAL-SUM)

Proof. (This proof remains unchanged as it concerns the estimators, not the update rule.) \square

G.2 PROOF OF PROPOSITION 2 (PER-PLAYER EXTERNAL REGRET WITH OMWU)

Proof. The proof adapts the standard analysis for Optimistic MWU to our setting with noisy, unbiased feedback for an arbitrary player p . The proof proceeds in two parts.

Part 1: Regret Against Estimated Losses. Let the estimated loss for player p be $\hat{l}_{t,i}^{(p)} = 1 - \hat{v}_{t,i}^{(p)}$. The OMWU algorithm uses an optimistic loss estimate $m_t^{(p)} = \hat{l}_t^{(p)} + (\hat{l}_t^{(p)} - \hat{l}_{t-1}^{(p)})$, with $\hat{l}_0^{(p)} = \mathbf{0}$. The weight update for player p is $w_{t+1}^{(p)}(i) = w_t^{(p)}(i) \exp(-\eta^{(p)} m_{t,i}^{(p)})$.

Let $W_t^{(p)} = \sum_{i=1}^{k_p(t)} w_t^{(p)}(i)$ be the potential function for player p . Following a standard potential function analysis (as detailed in Appendix C.1), the regret of OMWU with respect to the *observed* sequence of losses is bounded by its variation:

$$\sum_{t=1}^T \langle \sigma_t^{(p)}, \hat{l}_t^{(p)} \rangle - \min_i \sum_{t=1}^T \hat{l}_{t,i}^{(p)} \leq O \left(\sqrt{\ln k_{p,T} \sum_{t=1}^T \|\hat{l}_t^{(p)} - \hat{l}_{t-1}^{(p)}\|_\infty^2} \right) \quad (59)$$

Part 2: Translating to True Regret. We now relate this bound to the regret against the true payoffs $v_t^{(p)}$. Let the true loss be $l_t^{(p)} = 1 - v_t^{(p)}$ and the estimation error be $\Delta_t^{(p)} = \hat{l}_t^{(p)} - l_t^{(p)} = v_t^{(p)} - \hat{v}_t^{(p)}$. The cumulative true

regret for player p , $R_T^{(p)}$, is:

$$\begin{aligned}
R_T^{(p)} &= \sum_{t=1}^T \left(\max_i v_{t,i}^{(p)} - \langle \sigma_t^{(p)}, v_t^{(p)} \rangle \right) = \sum_{t=1}^T \left(\langle \sigma_t^{(p)}, l_t^{(p)} \rangle - \min_i l_{t,i}^{(p)} \right) \\
&= \sum_{t=1}^T \left(\langle \sigma_t^{(p)}, \hat{l}_t^{(p)} - \Delta_t^{(p)} \rangle - (\hat{l}_{t,i^*}^{(p)} - \Delta_{t,i^*}^{(p)}) \right) \\
&= \underbrace{\sum_{t=1}^T (\langle \sigma_t^{(p)}, \hat{l}_t^{(p)} \rangle - \hat{l}_{t,i^*}^{(p)})}_{\text{Regret on Estimates}} - \underbrace{\sum_{t=1}^T (\langle \sigma_t^{(p)}, \Delta_t^{(p)} \rangle - \Delta_{t,i^*}^{(p)})}_{\text{Cumulative Error Term}}
\end{aligned} \tag{60}$$

The first term is bounded as derived in Part 1. For the variation term in that bound, we use the triangle inequality:

$$\|\hat{l}_t^{(p)} - \hat{l}_{t-1}^{(p)}\|_\infty \leq \|\hat{l}_t^{(p)} - l_t^{(p)}\|_\infty + \|l_t^{(p)} - l_{t-1}^{(p)}\|_\infty + \|l_{t-1}^{(p)} - \hat{l}_{t-1}^{(p)}\|_\infty \tag{61}$$

$$\implies \|\hat{v}_t^{(p)} - \hat{v}_{t-1}^{(p)}\|_\infty \leq \|\Delta_t^{(p)}\|_\infty + \|v_t^{(p)} - v_{t-1}^{(p)}\|_\infty + \|\Delta_{t-1}^{(p)}\|_\infty \tag{62}$$

Taking the expectation over the sampling noise in our estimators, and noting that $\mathbb{E}[\Delta_t^{(p)}] = \mathbf{0}$ due to unbiasedness, we can bound the expected true regret. The error terms accumulate, leading to the final result. Dividing by T gives the statement in Proposition 56. \square

G.3 PROOF OF THEOREM 3 (PER-PLAYER ORACLE REGRET)

Proof. (This proof remains unchanged as it concerns the bandit oracle, which is decoupled from the meta-game solver by the two-time-scale assumption.) \square

G.4 PROOF OF THEOREM 4 (ϵ -CCE)

Proof. A time-averaged joint strategy $\bar{\mu}_T$ is an ϵ -CCE if for each player p and any deviating strategy π' , the gain from deviating is small:

$$\sum_{i,j} \bar{\mu}_T(i,j) r^{(p)}(\pi_i^{(1)}, \pi_j^{(2)}) \geq \sum_{i,j} \bar{\mu}_T(i,j) r^{(p)}(\pi', \pi_j^{(-p)}) - \epsilon_p \tag{63}$$

where $\epsilon = \sum_p \epsilon_p$. The maximum gain for player p from deviating is bounded by their average external regret against the sequence of opponent strategies $\{\sigma_t^{(-p)}\}_{t=1}^T$.

$$\epsilon_p \leq \frac{1}{T} \sum_{t=1}^T \left(\max_{\pi'} \mathbb{E}_{j \sim \sigma_t^{(-p)}} [r^{(p)}(\pi', \pi_j^{(-p)})] - \mathbb{E}_{i \sim \sigma_t^{(p)}, j \sim \sigma_t^{(-p)}} [r^{(p)}(\pi_i^{(p)}, \pi_j^{(-p)})] \right) \tag{64}$$

This is the definition of player p 's average exploitability of the sequence of joint strategies. We decompose this term for each player p , analogous to the decomposition in the proof of Theorem 3.4. The exploitability for player p at iteration t is:

$$\text{Exploit}_t^{(p)} = \underbrace{\left(\max_{i \in [k_p(t)]} (v_t^{(p)})_i - \langle \sigma_t^{(p)}, v_t^{(p)} \rangle \right)}_{\text{Internal Exploitability}^{(p)}} + \underbrace{\left(\max_{\pi'} \mathbb{E}_{j \sim \sigma_t^{(-p)}} [r^{(p)}(\pi', \pi_j^{(-p)})] - \max_{i \in [k_p(t)]} (v_t^{(p)})_i \right)}_{\text{Population Gap}^{(p)}} \tag{65}$$

Averaging over T and summing over players gives the total bound on ϵ .

1. **Internal Exploitability:** For each player p , this is their external regret within the restricted game. From Proposition 56, its average is bounded by the sum of the OMWU Regret and the MC Estimation Error.
2. **Population Gap:** For each player p , this term is decomposed further into their Amortized BR Error and Oracle Regret, following the same logic as in the proof of Theorem 3.4.

Summing these four distinct error sources for each player $p \in \{1, 2\}$ provides the composite bound for ϵ as stated in Theorem 58. \square

1276 H EXTENSION TO N-PLAYER GENERAL-SUM GAMES

1277
1278 We now generalize the framework to $n \geq 2$ players. Let $\mathcal{P} = \{1, \dots, n\}$ be the set of players, and for any player
1279 $p \in \mathcal{P}$, let $-p := \mathcal{P} \setminus \{p\}$ denote the set of all other players. At each iteration t , every player p maintains an
1280 anchor set $Z_t^{(p)}$, a generator G_{θ_p} , and a meta-strategy $\sigma_t^{(p)} \in \Delta_{k_p(t)-1}$. The joint mixture over all players is
1281 $\Sigma_t = \bigotimes_{q \in \mathcal{P}} \sigma_t^{(q)}$.
1282

1283 For each player p , the conceptual payoff is represented by a hypermatrix $M^{(p)}$ with entries corresponding to the
1284 expected reward for player p given a joint profile of pure strategies. The expected payoff for player p 's i -th policy
1285 against the joint mixture of all other players is:
1286

$$1287 v_t^{(p)}(i) = \mathbb{E}_{i_{-p} \sim \bigotimes_{q \in -p} \sigma_t^{(q)}} \left[r^{(p)}(\pi_i^{(p)}, \{\pi_{i_q}^{(q)}\}_{q \in -p}) \right] = \sum_{i_{-p}} M_{i, i_{-p}}^{(p)} \prod_{q \in -p} \sigma_t^{(q)}(i_q) \quad (66)$$

1288 and player p 's expected value under the full joint mixture is $\bar{r}_t^{(p)} = \sum_i \sigma_t^{(p)}(i) v_t^{(p)}(i)$.
1289

1292 H.1 MONTE CARLO META-GAME WITH IMPORTANCE WEIGHTING

1293 To avoid forming the computationally intractable payoff hypermatrices, we use importance-weighted estimators
1294 derived from a single set of shared game rollouts. We draw B joint strategy profiles $i^{(b)} = (i_b^{(1)}, \dots, i_b^{(n)}) \sim \Sigma_t$,
1295 run m episodes for each profile, and obtain returns $Y_{b,l}^{(p)}$ for every player p . The estimators for the per-policy value
1296 vector and the mixture value are:
1297

$$1298 \hat{v}_{t,i}^{(p)} = \frac{1}{Bm} \sum_{b=1}^B \sum_{l=1}^m \frac{\mathbf{1}\{i_b^{(p)} = i\}}{\sigma_t^{(p)}(i)} Y_{b,l}^{(p)} \quad (67)$$

$$1299 \hat{\bar{r}}_t^{(p)} = \frac{1}{Bm} \sum_{b=1}^B \sum_{l=1}^m Y_{b,l}^{(p)}$$

1300
1301 **Lemma H.1** (Unbiasedness and Concentration for n Players). *For rewards in $[0, 1]$, the estimators are unbiased:*
1302 $\mathbb{E}[\hat{v}_{t,i}^{(p)}] = v_t^{(p)}(i)$ and $\mathbb{E}[\hat{\bar{r}}_t^{(p)}] = \bar{r}_t^{(p)}$ for all p, i . *For any $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*
1303 *concentration bounds. However, the variance of the importance-weighted estimator $\hat{v}_{t,i}^{(p)}$ can be high if any policy's*
1304 *probability $\sigma_t^{(p)}(i)$ is small.*
1305

1312 H.2 PER-PLAYER OPTIMISTIC REPLICATOR AND ORACLE

1313 Each player p runs an independent learning process to update their meta-strategy and expand their policy set.

- 1314 • **OMWU Update:** Each player uses the Optimistic MWU rule to update their meta-strategy based on their
1315 individual payoff estimates. The optimistic estimate for player p is $m_t^{(p)} = 2\hat{v}_{t,i}^{(p)} - \hat{v}_{t-1}^{(p)}$, with $\hat{v}_0^{(p)} = \mathbf{0}$.

$$1316 \sigma_{t+1}^{(p)}(i) \propto \sigma_t^{(p)}(i) \exp \left(\eta_t^{(p)} [2\hat{v}_{t,i}^{(p)} - \hat{v}_{t-1,i}^{(p)} - \hat{\bar{r}}_t^{(p)}] \right) \quad (68)$$

- 1317 • **EB-UCB Oracle:** Each player p solves a separate multi-armed bandit problem to find a promising new
1318 anchor $z_t^{*,(p)}$. The value of an arm z is its expected reward against the joint mixture of all other players:

$$1319 f_t^{(p)}(z) = \mathbb{E}_{i_{-p} \sim \bigotimes_{q \in -p} \sigma_t^{(q)}} \left[r^{(p)}(\pi_{G_{\theta_p}}^{(p)}(z), \{\pi_{i_q}^{(q)}\}_{q \in -p}) \right] \quad (69)$$

1320
1321 **Proposition H.2** (Per-Player External Regret with OMWU in n-Player Games). *For any player p , assume payoffs*
1322 *lie in $[0, 1]$. The average external regret for player p using OMWU against the sequence of opponents' joint*
1323 *strategies is bounded by:*
1324

$$1325 \frac{1}{T} \sum_{t=1}^T \left(\max_i (v_t^{(p)})_i - \sigma_t^{(p)\top} v_t^{(p)} \right) \leq O \left(\frac{1}{T} \sqrt{\ln k_{p,T} \sum_{t=1}^T \|v_t^{(p)} - v_{t-1}^{(p)}\|_2^2} \right) \quad (70)$$

$$1326 + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\|\hat{v}_t^{(p)} - v_t^{(p)}\|_\infty \right]$$

Theorem H.3 (Per-Player Instance-Dependent Oracle Regret). *Under a two-time-scale assumption for n -players, where opponents' mixtures $\{\sigma_t^{(q)}\}_{q \in -p}$ are considered fixed during player p 's oracle step, the cumulative oracle regret for each player is bounded:*

$$\sum_{t=1}^T \mathbb{E}[f_t^{(p)}(z_t^{*,(p)}) - f_t^{(p)}(z_t^{(p)})] = O\left(\sum_{z \neq z^{*,(p)}} \frac{\ln T}{\Delta_z^{(p)}}\right) + \lambda_{J,p} \sum_{t=1}^T \mathbb{E}\left[\|J_{G_{\theta_p}}(z_t^{*,(p)})\|_F^2\right] \quad (71)$$

H.3 OVERALL GUARANTEE: ϵ -CCE FOR N PLAYERS

The time-averaged joint play, defined as $\bar{\mu}_T(i_1, \dots, i_n) := \frac{1}{T} \sum_{t=1}^T \prod_{p=1}^n \sigma_t^{(p)}(i_p)$, converges to an ϵ -Coarse-Correlated Equilibrium (ϵ -CCE).

Theorem H.4 (ϵ -CCE of Time-Average Joint Play for n Players). *Assume rewards in $[0, 1]$ and the n -player two-time-scale assumption. With per-player OMWU meta-solvers, the time-averaged distribution $\bar{\mu}_T$ is an ϵ -CCE, where the total deviation incentive ϵ is bounded by the sum of regrets and errors across all players:*

$$\begin{aligned} \epsilon \leq \sum_{p=1}^n \left[\underbrace{O\left(\frac{1}{T} \sqrt{\ln k_{p,T} \sum_{t=1}^T \|v_t^{(p)} - v_{t-1}^{(p)}\|_\infty^2}\right)}_{\text{OMWU Regret}} + \underbrace{\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t^{(p)} - v_t^{(p)}\|_\infty]}_{\text{MC Estimation Error}} \right. \\ \left. + \underbrace{\frac{1}{T} \sum_{t=1}^T \epsilon_{\text{BR},t}^{(p)}}_{\text{Amortized BR Error}} + \underbrace{O\left(\frac{1}{T} \sum_{z \neq z^{*,(p)}} \frac{\ln T}{\Delta_z^{(p)}}\right)}_{\text{Oracle Regret}} \right] \quad (72) \end{aligned}$$

If simulation and training budgets grow such that the error terms for each player go to zero, then $\epsilon \rightarrow 0$ as $T \rightarrow \infty$.

I PROOFS FOR N-PLAYER GENERAL-SUM

I.1 PROOF OF LEMMA 3 (UNBIASEDNESS AND CONCENTRATION, N-PLAYERS)

Proof. (This proof remains unchanged as it concerns the estimators, not the update rule.) \square

I.2 PROOF OF PROPOSITION 3 (PER-PLAYER EXTERNAL REGRET WITH OMWU)

Proof. The proof is structurally identical to the proof of Proposition 2 for the two-player case. We present it here explicitly for the n -player setting for completeness. The analysis is performed for an arbitrary player $p \in \{1, \dots, n\}$.

Part 1: Regret Against Estimated Losses. Let the estimated loss for player p be $\hat{l}_{t,i}^{(p)} = 1 - \hat{v}_{t,i}^{(p)}$. The OMWU algorithm uses an optimistic loss estimate $m_t^{(p)} = \hat{l}_t^{(p)} + (\hat{l}_t^{(p)} - \hat{l}_{t-1}^{(p)})$, with $\hat{l}_0^{(p)} = \mathbf{0}$. The weight update for player p is $w_{t+1}^{(p)}(i) = w_t^{(p)}(i) \exp(-\eta^{(p)} m_{t,i}^{(p)})$.

The regret of OMWU with respect to this observed sequence of losses is bounded by its variation, following a standard potential function analysis:

$$\sum_{t=1}^T \langle \sigma_t^{(p)}, \hat{l}_t^{(p)} \rangle - \min_i \sum_{t=1}^T \hat{l}_{t,i}^{(p)} \leq O\left(\sqrt{\ln k_{p,T} \sum_{t=1}^T \|\hat{l}_t^{(p)} - \hat{l}_{t-1}^{(p)}\|_\infty^2}\right) \quad (73)$$

The "environment" from player p 's perspective is the sequence of payoff vectors it receives. The regret bound does not depend on how this sequence is generated by the other $n - 1$ players, only on its properties (i.e., its variation).

Part 2: Translating to True Regret. We relate this bound to the regret against the true payoffs $v_t^{(p)}$. Let the true loss be $l_t^{(p)} = 1 - v_t^{(p)}$ and the estimation error be $\Delta_t^{(p)} = \hat{l}_t^{(p)} - l_t^{(p)} = v_t^{(p)} - \hat{v}_t^{(p)}$. The cumulative true regret

for player p , $R_T^{(p)}$, is:

$$\begin{aligned}
R_T^{(p)} &= \sum_{t=1}^T \left(\max_i v_{t,i}^{(p)} - \langle \sigma_t^{(p)}, v_t^{(p)} \rangle \right) = \sum_{t=1}^T \left(\langle \sigma_t^{(p)}, l_t^{(p)} \rangle - \min_i l_{t,i}^{(p)} \right) \\
&= \underbrace{\sum_{t=1}^T \left(\langle \sigma_t^{(p)}, \hat{l}_t^{(p)} \rangle - \min_i \hat{l}_{t,i}^{(p)} \right)}_{\text{Regret on Estimates}} - \underbrace{\sum_{t=1}^T \left(\langle \sigma_t^{(p)}, \Delta_t^{(p)} \rangle - \Delta_{t,i^*}^{(p)} \right)}_{\text{Cumulative Error Term}}
\end{aligned} \tag{74}$$

The first term is bounded as derived in Part 1. The variation term in that bound is handled with the triangle inequality: $\|\hat{v}_t^{(p)} - \hat{v}_{t-1}^{(p)}\|_\infty \leq \|\Delta_t^{(p)}\|_\infty + \|v_t^{(p)} - v_{t-1}^{(p)}\|_\infty + \|\Delta_{t-1}^{(p)}\|_\infty$. Taking the expectation over the sampling noise and dividing by T yields the final bound stated in the proposition. \square

I.3 PROOF OF THEOREM 5 (PER-PLAYER ORACLE REGRET)

Proof. The proof is a standard bandit analysis, applied to the independent learning problem faced by each player p . Under the two-time-scale assumption, the joint mixture of all other players, $\Sigma_t^{(-p)} = \otimes_{q \in -p} \sigma_t^{(q)}$, is considered fixed during player p 's oracle selection phase at meta-iteration t .

For player p , the learning problem is a stochastic multi-armed bandit task where:

- The set of arms is the candidate pool of latent codes, $\Lambda_t^{(p)}$.
- The reward for pulling an arm $z \in \Lambda_t^{(p)}$ is a random variable (the outcome of a game rollout) whose expectation is the true arm value, $f_t^{(p)}(z)$, as defined in Equation 69.

This is a standard bandit setting. The EB-UCB algorithm guarantees that the number of times a suboptimal arm z is pulled, $N_z(T)$, is bounded logarithmically with respect to the total number of oracle steps T , i.e., $\mathbb{E}[N_z(T)] = O(\ln T / \Delta_z^{(p)})$, where $\Delta_z^{(p)}$ is the suboptimality gap. The total regret is the sum of the expected regrets from pulling each suboptimal arm. The Jacobian penalty adds a simple additive term to this regret, leading to the final bound stated in the main text. \square

I.4 PROOF OF THEOREM 6 (ϵ -CCE FOR N PLAYERS)

Proof. The proof generalizes the argument from the two-player case. A time-averaged joint distribution $\bar{\mu}_T$ is an ϵ -CCE if for every player $p \in \{1, \dots, n\}$ and for every possible unilateral deviation to a pure strategy π'_p , the incentive to deviate is bounded:

$$\mathbb{E}_{i \sim \bar{\mu}_T} [r^{(p)}(\pi_{i_p}^{(p)}, \pi_{i_{-p}}^{(-p)})] \geq \mathbb{E}_{i_{-p} \sim \bar{\mu}_T^{(-p)}} [r^{(p)}(\pi'_p, \pi_{i_{-p}}^{(-p)})] - \epsilon_p \tag{75}$$

where $\epsilon = \sum_{p=1}^n \epsilon_p$. The maximum gain for player p from deviating is bounded by their average external regret over the T iterations. Let's bound this gain, ϵ_p :

$$\epsilon_p \leq \frac{1}{T} \sum_{t=1}^T \left(\max_{\pi'_t} \mathbb{E}_{i_{-p} \sim \Sigma_t^{(-p)}} [r^{(p)}(\pi'_t, \pi_{i_{-p}}^{(-p)})] - \mathbb{E}_{i \sim \Sigma_t} [r^{(p)}(\pi_i)] \right) \tag{76}$$

This is precisely player p 's average exploitability over the sequence of play. We decompose this term for each player p into four components, following the structure of the proof of Theorem 3.4. For each player p , their average exploitability is bounded by the sum of:

1. **Average Internal Exploitability:** The regret of OMWU within the restricted game on $Z_t^{(p)}$, which is bounded by Proposition 3. This gives the **OMWU Regret** and the **MC Estimation Error** terms.
2. **Average Population Gap:** The gap between a true best response and the best response within $Z_t^{(p)}$. This is further decomposed into the **Amortized BR Error** (by definition) and the **Oracle Regret** (from Theorem 5).

Thus, for each player p , their maximum deviation incentive ϵ_p is bounded by the sum of their four error terms:

$$\epsilon_p \leq O \left(\frac{1}{T} \sqrt{\ln k_{p,T} \sum_{t=1}^T \|v_t^{(p)} - v_{t-1}^{(p)}\|_\infty^2} \right) + \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\hat{v}_t^{(p)} - v_t^{(p)}\|_\infty] + \bar{\epsilon}_{\text{BR}}^{(p)} + \bar{\mathcal{R}}_T^{(p)} \tag{77}$$

1450 where $\bar{\varepsilon}_{\text{BR}}^{(p)}$ and $\bar{\mathcal{R}}_T^{(p)}$ are the average BR error and oracle regret, respectively. The total error ϵ is the sum of these
1451 deviation incentives, $\epsilon = \sum_{p=1}^n \epsilon_p$. Summing the bounds for each player provides the composite bound for ϵ as
1452 stated in Theorem 72. \square
1453
1454

1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507

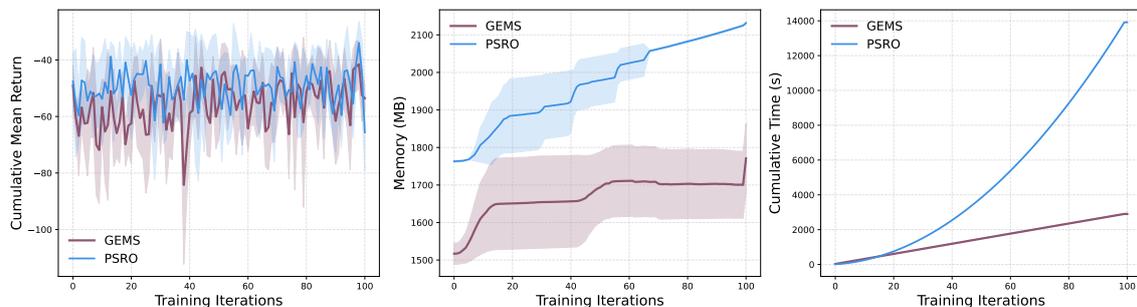
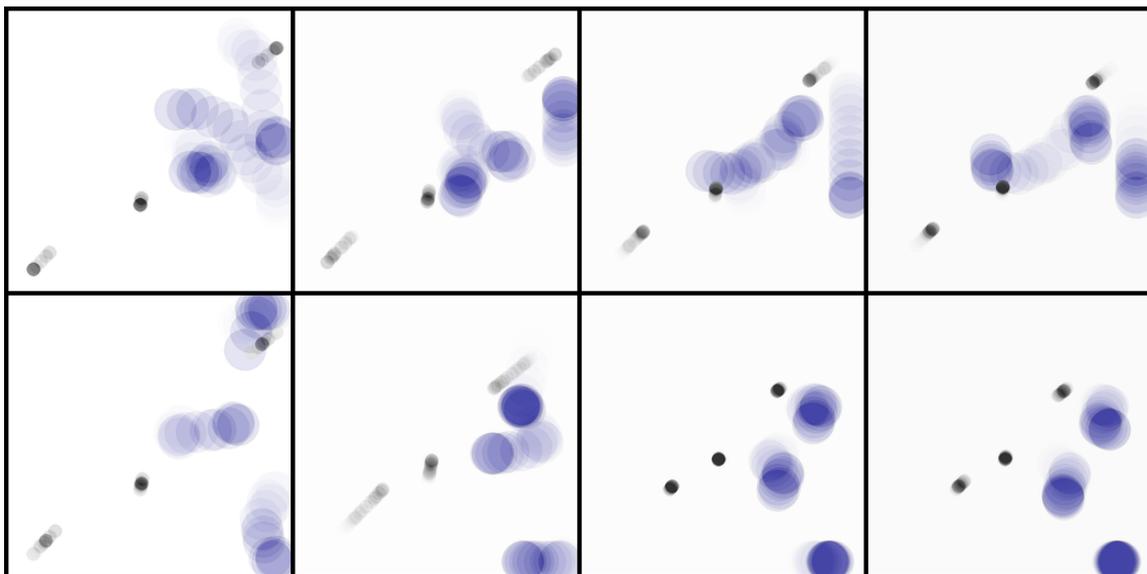
Part III

Experiments

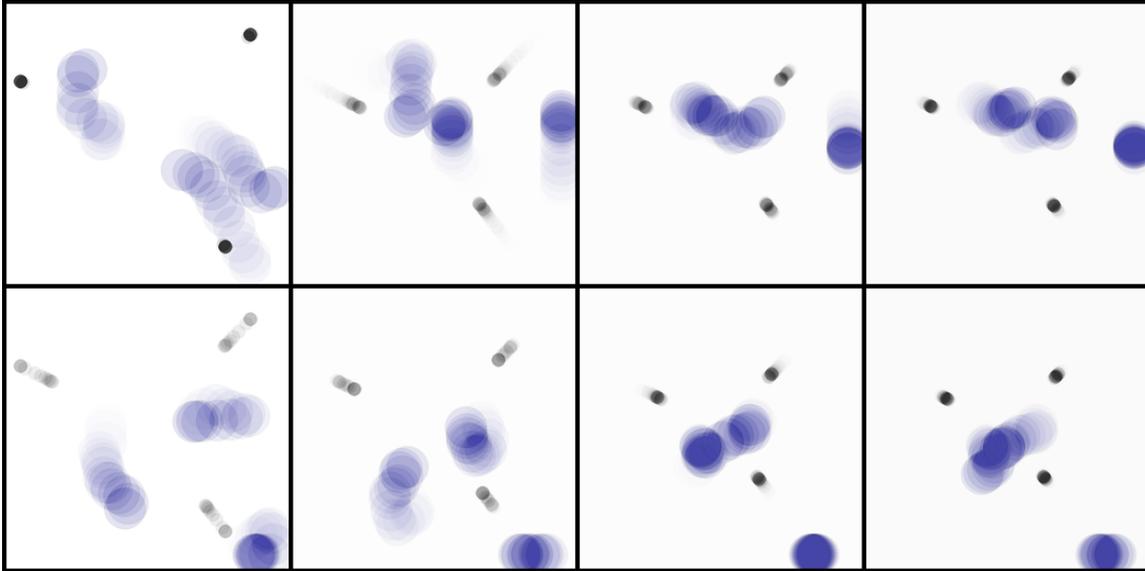
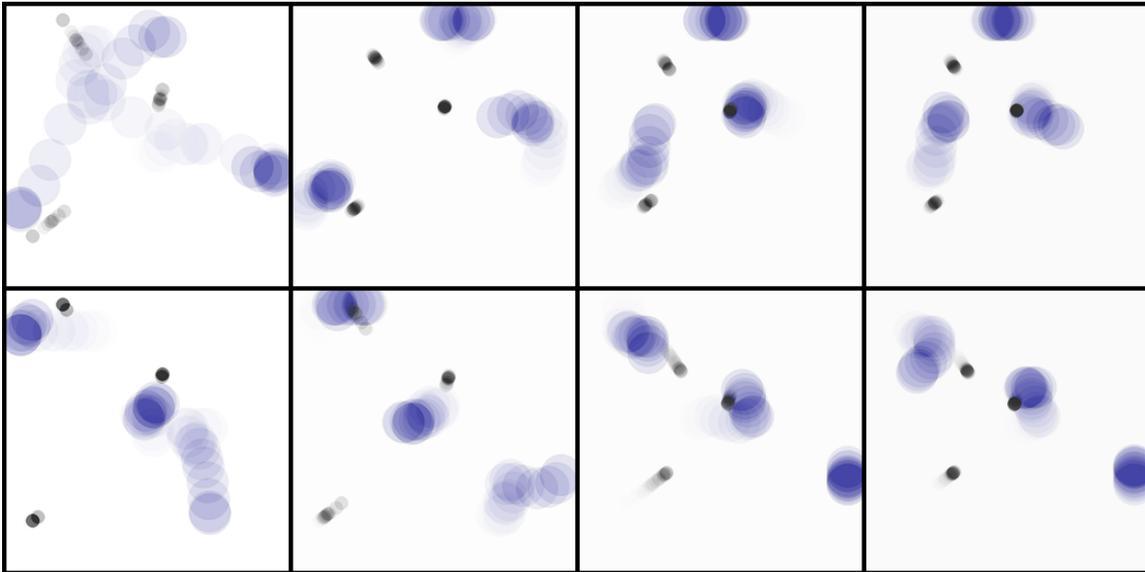
J COORDINATION ON SIMPLE SPREAD

Setup. Our second experiment is conducted in the `simple_spread` environment, a classic cooperative benchmark from PettingZoo. The task requires three agents to collaboratively “cover” three distinct target landmarks in a 2D space. Agents are rewarded for minimizing their distance to any landmark but are penalized for colliding with each other. A successful outcome requires the agents to learn a decentralized “divide and conquer” strategy, assigning themselves to unique targets and navigating to them efficiently. This environment is designed to test emergent cooperation and task division. We again compare GEMS and PSRO on qualitative behavior, mean return, memory, and computation time, averaged over 5 seeds.

Objective. This experiment aims to achieve two goals. First, to confirm the significant scalability benefits of GEMS, as demonstrated in the previous experiment, but now in a purely cooperative setting. Second, to evaluate which algorithm is more effective at discovering the complex, coordinated strategies required for cooperative success, as measured by both quantitative rewards and qualitative analysis of agent behaviors.

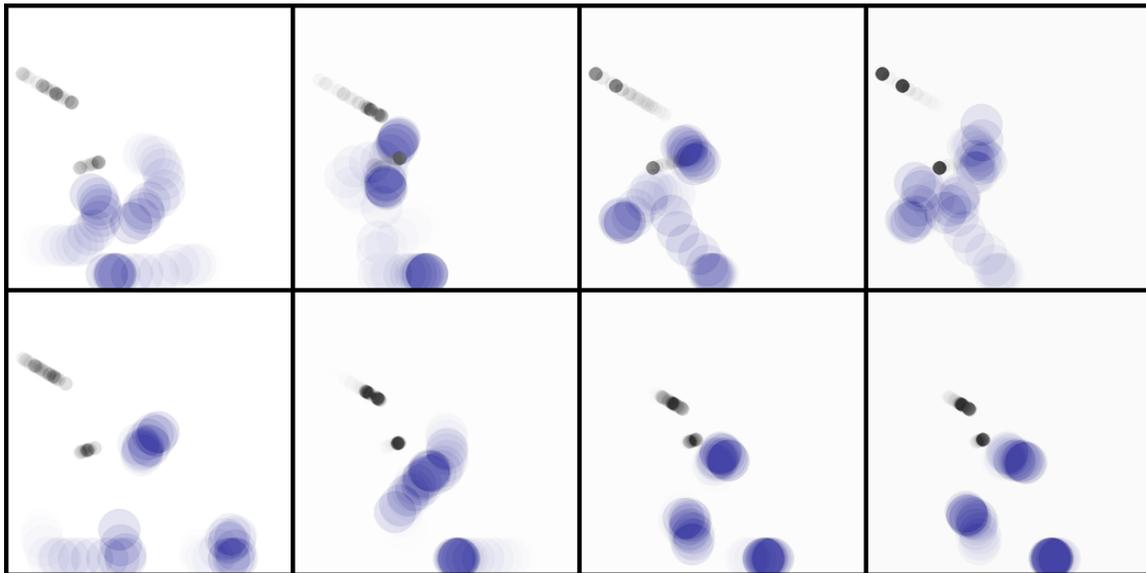
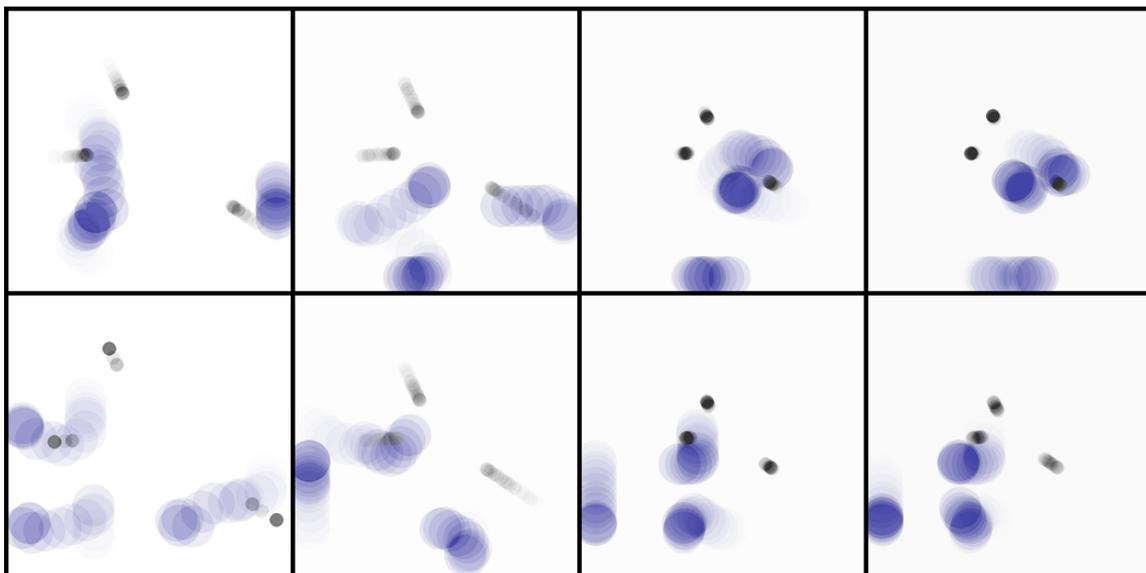
Figure 6: *Simple Spread*Figure 7: *Simple Spread* run on *Seed 0*

Results. Across `simple_spread`, GEMS again outperforms PSRO in both effectiveness and efficiency. In the ghost-montage comparisons for seeds 0–4 (Figures 7–11; GEMS on top, PSRO on bottom), **GEMS agents learn an**

Figure 8: Simple Spread run on *Seed 1*Figure 9: Simple Spread run on *Seed 2*

efficient, coordinated strategy: they quickly partition landmarks and move directly toward them, and when two agents collide they *separate and re-plan* rather than dithering—while still staying close to the black-dot targets (agents rendered in blue). By contrast, **PSRO agents exhibit weaker coordination:** typically one–two agents hover near a target while another drifts away from objectives or circles indecisively, an undesirable failure mode visible in multiple seeds (e.g., Fig. 7, 8). This qualitative advantage aligns with the aggregate plot (Figure 6): GEMS achieves a higher (less negative) mean return with lower variance over time. In terms of scalability, results match prior sections: GEMS is over 6× faster in cumulative time (~2,000s vs. ~13,000s) and maintains a flat memory profile, while PSRO’s resource usage scales poorly.

Analysis. These rollouts (Figures 7–11) reinforce our central claim: GEMS is better suited for cooperative multi-agent tasks that require implicit role assignment. The **EB-UCB oracle’s exploration over a diverse latent space** helps discover complementary roles (distinct landmark assignments) and preserves them under contact—agents momentarily repel when they touch, then settle back near their respective black-dot objectives. Simpler procedures in PSRO more often collapse to near-symmetric yet suboptimal behaviors (e.g., multiple agents chasing the same landmark while another disengages). Coupled with GEMS’s **surrogate-free design**—which yields lower wall-clock time and stable memory—these qualitative and quantitative results position GEMS as a robust, scalable framework for cooperative MARL.

Figure 10: Simple Spread run on *Seed 3*Figure 11: Simple Spread run on *Seed 4*

K COORDINATION ON SIMPLE TAG

Simple Tag (Appendix). Across seeds 0–4 (ghost montages: GEMS on top, PSRO on bottom), **GEMS** reliably produces *coordinated pursuit*: the red taggers self-organize into a roughly triangular encirclement that narrows angles on the green runner, adjusting spacing to cut off escape lanes; the green agent, in turn, exhibits purposeful evasion (arcing and zig-zag trajectories), and *captures do occur* when the enclosure closes. Under **PSRO**, by contrast, the taggers seldom sustain a stable formation: transient “triangle-ish” shapes appear but collapse before containment, leaving wide gaps through which the green agent *easily escapes*. While the green agent under PSRO shows less consistently trained avoidance than under GEMS, the taggers’ lack of persistent coordination dominates the outcome, captures are rare and evasion is commonplace. Overall, these qualitative rollouts corroborate the main-paper claim: GEMS induces higher-level cooperative tactics (persistent enclosure and angle closing) that the PSRO baseline fails to discover or maintain.

1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739

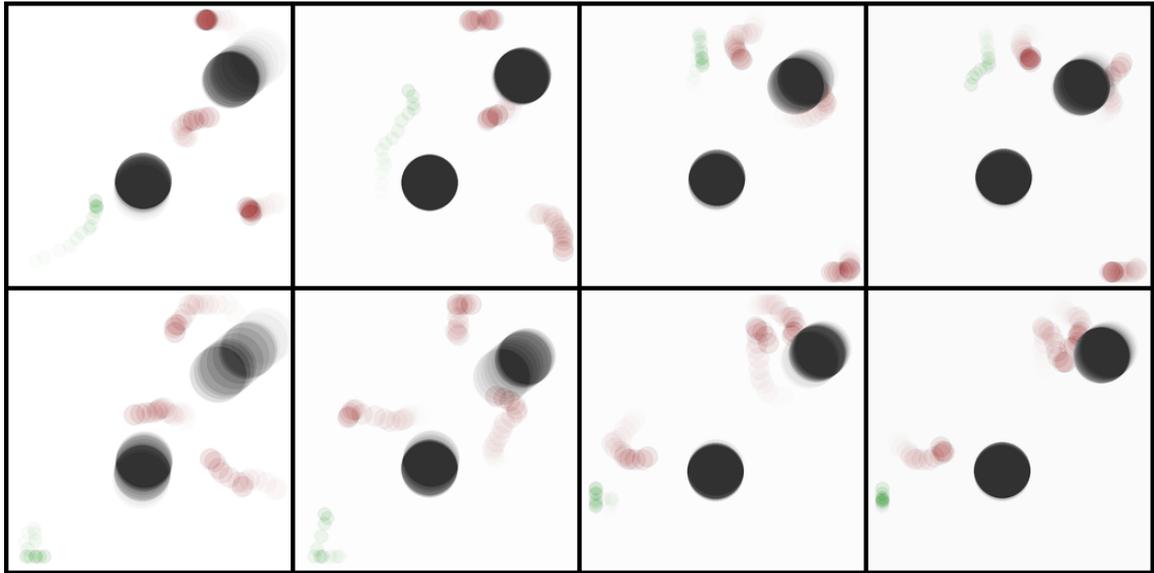


Figure 12: Simple Tag run on *Seed 0*

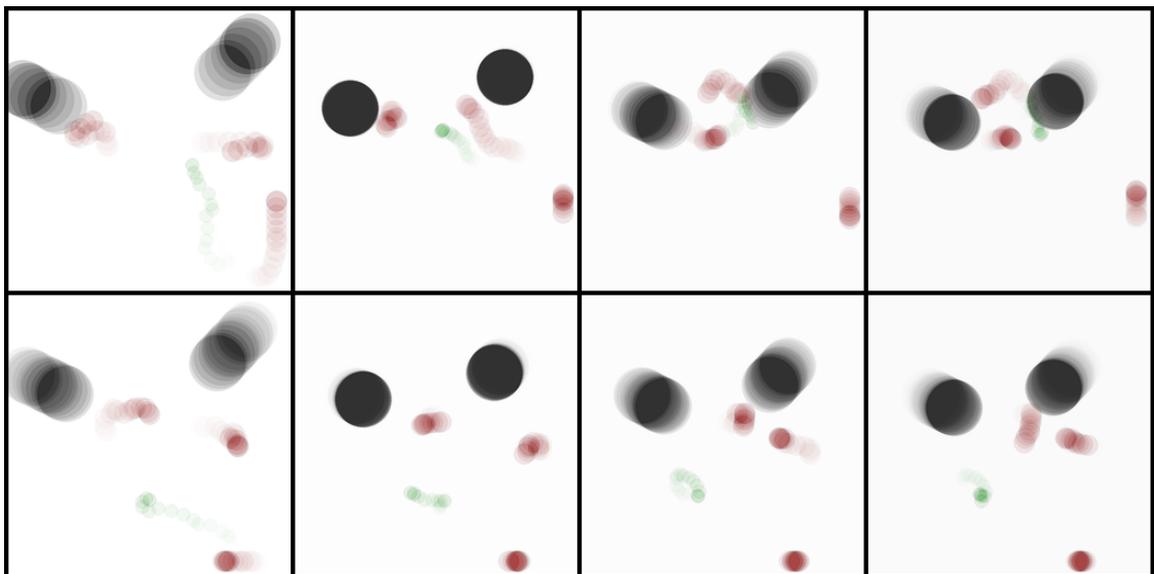


Figure 13: Simple Tag run on *Seed 1*

1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797

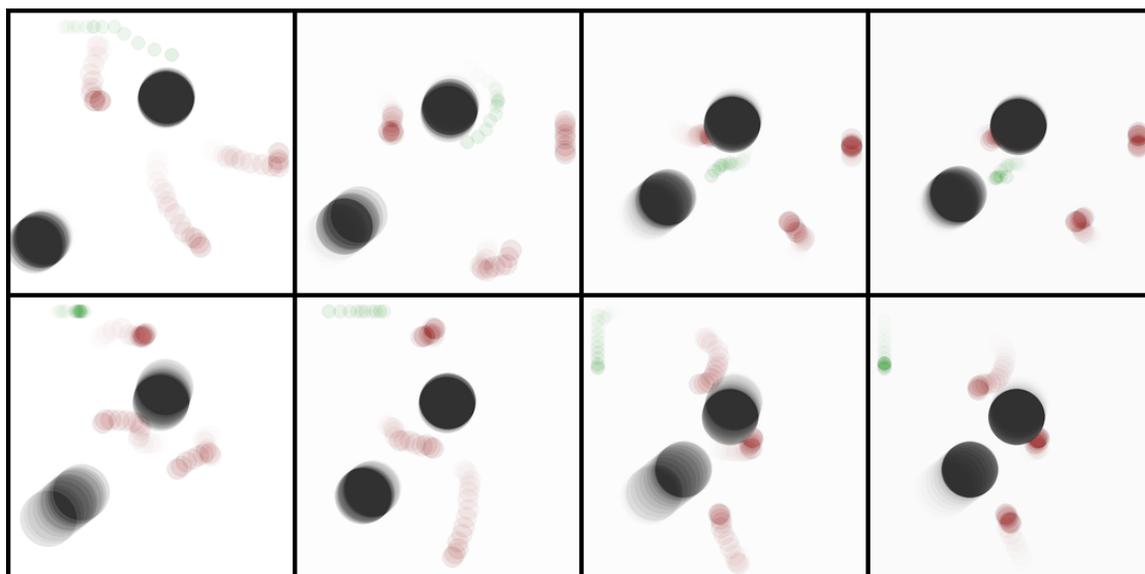


Figure 14: Simple Tag run on *Seed 2*

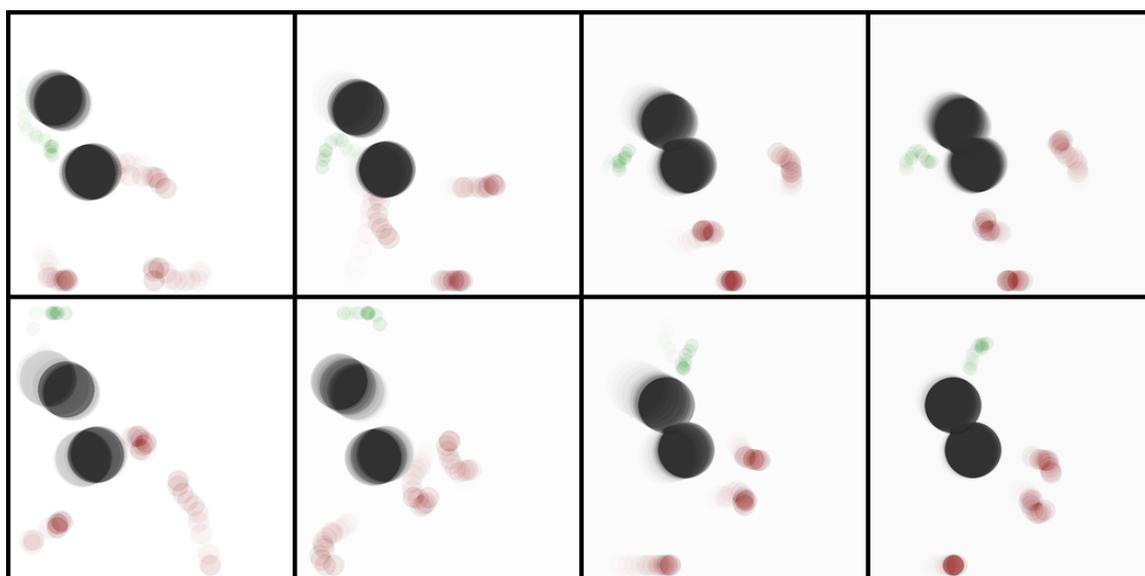


Figure 15: Simple Tag run on *Seed 3*

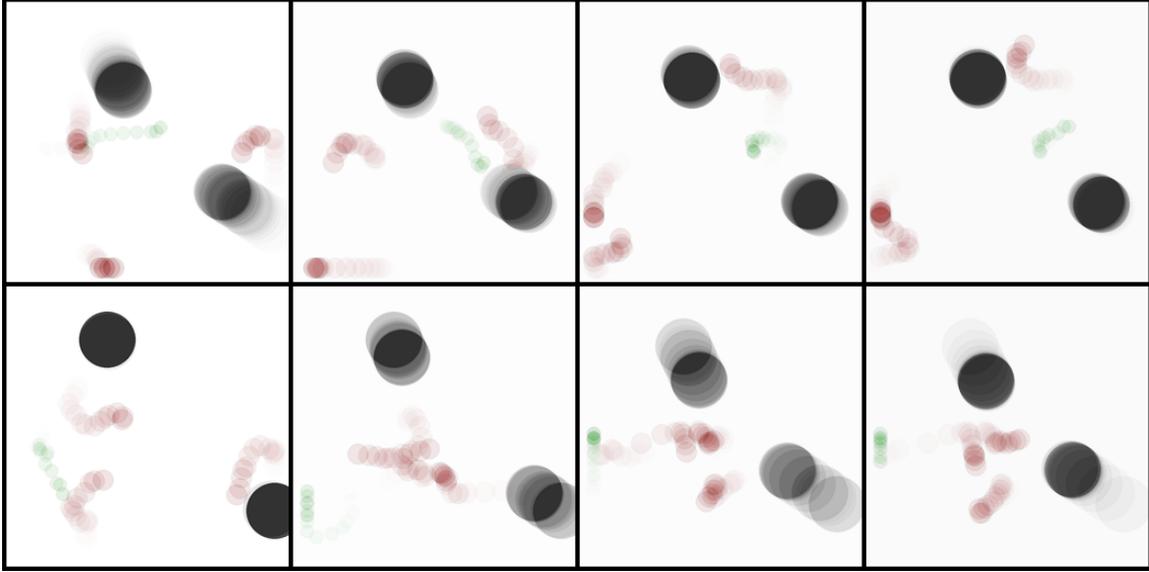


Figure 16: Simple Tag run on Seed 4

L ABLATION ON PUBLIC GOODS GAME

Public Goods Game Ablation: Purpose. This ablation studies *which parts of FAST GEMS-OMWU matter most* for outcome quality versus compute, so that our defaults are principled and robust rather than hand-tuned. Concretely, we vary algorithmic and stability knobs— η , η_{sched} (const/sqrt/harmonic), ema , mwu_grad_cap , temperature τ ; population-growth and exploration knobs— oracle_nz , oracle_period , pool_mut , pool_rand , latent size zdim ; and ABR-TR tightness— abr_lr , abr_steps , and β_{KL} . We summarize effects with terminal and speed-sensitive metrics: welfare_last , coop_last (final social welfare and cooperation) and their learning-curve integrals welfare_auc , coop_auc , alongside efficiency measures $\text{time_total_sec_avg}$ and ram_peak_mb_avg . Because a full factorial sweep is combinatorially infeasible, we run structured partial variations (one-factor and a few targeted pairs) across multiple seeds and report $\text{mean} \pm \text{std}$ from the multi-seed runner; the intent is to **characterize sensitivities and trade-offs**, not to claim a single global optimum. In practice, the ablation clarifies (i) step-size scheduling and mild smoothing (harmonic η_t with modest EMA) stabilize meta-updates without slowing progress unduly, (ii) moderate β_{KL} in ABR-TR improves welfare without large time/RAM penalties, and (iii) small increases in oracle_period and balanced $\text{pool_mut}/\text{pool_rand}$ control population growth while preserving exploration. All runs log the exact configuration to CSV for reproducibility and allow further probing of interactions as needed.

Public Goods Game: overview & metrics. The *Public Goods Game (PGG)* is an n -player social-dilemma benchmark. Each player chooses to *contribute* (cooperate) or *withhold* (defect). Let $a_i \in \{0, 1\}$ denote player i 's contribute decision, $S = \sum_j a_j$ the total contributions, multiplier $r > 0$, and per-contribution cost $c_i > 0$ (homogeneous case uses $c_i \equiv c$). The one-round payoff to player i is

$$u_i = \frac{r}{n} S - c_i a_i, \quad \text{and the social welfare } W = \sum_{i=1}^n u_i = r S - \sum_{i=1}^n c_i a_i. \quad (78)$$

When $\frac{r}{n} < c < r$ (homogeneous costs), each individual gains by defecting (free-riding) while the group gains by coordinating on cooperation—this is the core tension PGG exposes.

In our runs, the generator outputs the probability of cooperation; we track both *outcome quality* and *learning speed/stability* via four summary metrics:

- **welfare_last**: the mean social welfare W at the *final training iteration* (averaged over profile samples and seeds). Higher is better; if parameters make W negative, “less negative” still indicates improvement.
- **coop_last**: the final average cooperation rate $\frac{1}{n} \sum_i \Pr[a_i=1]$.
- **welfare_auc**: the (trapezoidal) area under the welfare-versus-iteration curve. This rewards methods that reach good welfare earlier and keep it stable—a joint measure of speed and robustness.

- **coop_auc**: the analogous AUC for cooperation rate.

All four are reported as mean±std across seeds by the multi-seed runner, making it straightforward to compare algorithms and hyper-parameter settings on both efficiency and cooperative outcomes.

iters_in_child = 10; seeds_in_child = 0,1,2,3,4; iters = 10; n_players = 5.

Table 1: Ablation (Public Goods Game)

welfare_last	coop_last	welfare_auc	coop_auc	time_total_sec_avg	ram_peak_mib_avg	abr_lr	abr_steps	beta_kl	cost_c	delta0	ema	eta	eta_sched	hetero	lambda_jac	log_ubc	mult_r	mwu_grad_cap	no_plots	num_seeds	oracle_nz	oracle_period	pool_mnt	pool_rand	slow	slowdown	tau	top_seed	ubc_nz	zdim
0.008177475	0.000817748	0.997798743	0.099779875	1.252423239	1240.809375	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	8	
0.0481612	0.00481612	4.389817474	0.438981745	1.243599987	1240.700781	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	8	
0.007937888	0.000793789	0.899287799	0.08992878	1.250558043	1240.736719	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	8	
0.036606287	0.003660629	4.423422511	0.442342253	1.220253325	1240.854688	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	8
0.021302772	0.002130277	1.155576756	0.115557676	1.270089245	1241.642188	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	16	
0.065991034	0.006599103	4.265299203	0.42652992	1.26079936	1240.81125	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	16	
0.009342633	0.000934263	0.683710212	0.068371021	1.226696777	1240.733594	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	16	
0.031212943	0.003121294	3.660477072	0.366047713	1.247665262	1240.683594	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	16
0.008177475	0.000817748	0.997798743	0.099779875	1.252423239	1240.809375	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	8	
0.0481612	0.00481612	4.389817474	0.438981745	1.243599987	1240.700781	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	8	
0.007937888	0.000793789	0.899287799	0.08992878	1.250558043	1240.736719	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	8	
0.036606287	0.003660629	4.423422511	0.442342253	1.220253325	1240.854688	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	8
0.021302772	0.002130277	1.155576756	0.115557676	1.270089245	1241.642188	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	16	
0.065991034	0.006599103	4.265299203	0.42652992	1.26079936	1240.81125	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	16	
0.009342633	0.000934263	0.683710212	0.068371021	1.226696777	1240.733594	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	16	
0.031212943	0.003121294	3.660477072	0.366047713	1.247665262	1240.683594	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	16
0.008177475	0.000817748	0.997798743	0.099779875	1.252423239	1240.809375	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	8	
0.0481612	0.00481612	4.389817474	0.438981745	1.243599987	1240.700781	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	8	
0.007937888	0.000793789	0.899287799	0.08992878	1.250558043	1240.736719	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	8	
0.036606287	0.003660629	4.423422511	0.442342253	1.220253325	1240.854688	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	8
0.021302772	0.002130277	1.155576756	0.115557676	1.270089245	1241.642188	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	16	
0.065991034	0.006599103	4.265299203	0.42652992	1.26079936	1240.81125	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	16	
0.009342633	0.000934263	0.683710212	0.068371021	1.226696777	1240.733594	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	16	
0.031212943	0.003121294	3.660477072	0.366047713	1.247665262	1240.683594	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	16
0.008177475	0.000817748	0.997798743	0.099779875	1.252423239	1240.809375	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	8	
0.0481612	0.00481612	4.389817474	0.438981745	1.243599987	1240.700781	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	8	
0.007937888	0.000793789	0.899287799	0.08992878	1.250558043	1240.736719	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	8	
0.036606287	0.003660629	4.423422511	0.442342253	1.220253325	1240.854688	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	8
0.021302772	0.002130277	1.155576756	0.115557676	1.270089245	1241.642188	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	16	
0.065991034	0.006599103	4.265299203	0.42652992	1.26079936	1240.81125	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	16	
0.009342633	0.000934263	0.683710212	0.068371021	1.226696777	1240.733594	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	16	
0.031212943	0.003121294	3.660477072	0.366047713	1.247665262	1240.683594	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	16
0.008177475	0.000817748	0.997798743	0.099779875	1.252423239	1240.809375	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	8	
0.0481612	0.00481612	4.389817474	0.438981745	1.243599987	1240.700781	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	8	
0.007937888	0.000793789	0.899287799	0.08992878	1.250558043	1240.736719	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	8	
0.036606287	0.003660629	4.423422511	0.442342253	1.220253325	1240.854688	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	8
0.021302772	0.002130277	1.155576756	0.115557676	1.270089245	1241.642188	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	16	
0.065991034	0.006599103	4.265299203	0.42652992	1.26079936	1240.81125	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	16	
0.009342633	0.000934263	0.683710212	0.068371021	1.226696777	1240.733594	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	16	
0.031212943	0.003121294	3.660477072	0.366047713	1.247665262	1240.683594	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	1	1	2	1.5	0	8	16
0.008177475	0.000817748	0.997798743	0.099779875	1.252423239	1240.809375	0.0005	20	0.01	1	0.0005	0.0	0.08	const	0	0	0	3	1	5	1	1	2	1	0	1	1.0	0	8	8	
0.0481612	0.00481612	4.389817474	0.438981745	1.243599987	1240.700781	0.00025	20	0.02	1	0.0005	0.0	0.04	const	0	0	0	3	1	5	1	1	2	1	0	2	1.5	0	8	8	
0.007937888	0.000793789	0.899287799	0.08992878	1.250558043	1240.736719	0.0005	20	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	1	1	1.0	0	8	8	
0.036606287	0.003660629	4.423422511	0.442342253	1.220253325	1240.854688	0.00025	20	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05												

1914
1915
1916
1917
1918

	welfare_last	coop_last	welfare_auc	coop_auc	time_total_sec_avg	ram_peak_mb_avg	abr_lr	abr_steps	beta_kl	cos_c	delta0	ema	eta	eta_sched	hetero	lambda_jac	log_ubc	mult_r	mwu_grad_cap	no_plots	num_seeds	oracle_nz	oracle_period	pool_mut	pool_rand	slow	slowdown	tau	top_seed	ubc_nz	zdim
1919	0.012002162	0.001200216	1.105425272	0.110542528	2.323050165	1240.639844	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	2	2	2	1	1	2	1.5	0	8	16
1920	0.006682226	0.000668223	0.315112181	0.031511218	2.385132074	1241.5875	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	1	1	2	2	0	1	1.0	0	8	8
1921	0.025486331	0.002548633	1.723243252	0.172324327	2.37188077	1240.865625	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	1	1	2	2	0	2	1.5	0	8	8
1922	0.000513562	5.14e-05	0.091702681	0.009170268	2.325033045	1240.753125	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	2	1	1	1.0	0	8	8
1923	0.008697222	0.000869722	1.276544507	0.127654448	2.342204618	1240.599219	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	2	1	2	1.5	0	8	8
1924	0.012916667	0.001291667	0.387960172	0.038796017	2.324785471	1240.727344	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	1	1	2	2	0	1	1.0	0	8	16
1925	0.031189462	0.003118946	1.559489052	0.155948906	2.366880798	1240.7125	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	1	1	2	2	0	2	1.5	0	8	16
1926	0.002724909	0.000272491	0.117010106	0.011701011	2.305494022	1240.871875	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	2	2	1	1	1.0	0	8	16
1927	0.019667632	0.001966763	1.234394693	0.12343947	2.265299177	1240.741406	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	2	2	1	2	1.5	0	8	16
1928	0.006682226	0.000668223	0.315112181	0.031511218	2.417394686	1240.590625	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	2	1	2	2	0	1	1.0	0	8	8
1929	0.025486331	0.002548633	1.723243252	0.172324327	2.331274128	1240.888281	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	2	1	2	2	0	2	1.5	0	8	8
1930	0.000513562	5.14e-05	0.091702681	0.009170268	2.322550297	1240.583594	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	2	2	2	2	1	1	1.0	0	8	8
1931	0.008697222	0.000869722	1.276544507	0.127654448	2.306919479	1240.950781	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	2	2	2	2	1	2	1.5	0	8	8
1932	0.012916667	0.001291667	0.387960172	0.038796017	2.33807683	1241.0	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	2	1	2	2	0	1	1.0	0	8	16
1933	0.031189462	0.003118946	1.559489052	0.155948906	2.388987923	1240.896094	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	2	1	2	2	0	2	1.5	0	8	16
1934	0.002724909	0.000272491	0.117010106	0.011701011	2.299126863	1240.898438	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	2	2	2	2	1	1	1.0	0	8	16
1935	0.019667632	0.001966763	1.234394693	0.12343947	2.294810247	1240.71875	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	2	2	2	2	1	2	1.5	0	8	16
1936	0.007332874	0.000733287	0.285362167	0.028536217	2.389500189	1241.516406	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	1	1	4	1	0	1	1.0	0	8	8
1937	0.03894474	0.003894474	1.836211315	0.18362113	2.413767529	1240.792188	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	1	1	4	1	0	2	1.5	0	8	8
1938	0.000533205	5.33e-05	0.087070911	0.008707091	2.328015566	1241.722656	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	4	1	1	1	1.0	0	8	8
1939	0.008562821	0.000856282	1.201094464	0.120109444	2.327522945	1240.875781	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	4	1	1	2	1.5	0	8	8
1940	0.004823769	0.000482377	0.222542943	0.022254294	2.311411142	1240.721094	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	1	1	4	1	0	1	1.0	0	8	16
1941	0.031591403	0.00315914	1.51600726	0.151600726	2.387298203	1240.904688	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	1	1	4	1	0	2	1.5	0	8	16
1942	0.001338041	0.000133804	0.084246678	0.008424668	2.302635479	1240.71875	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	4	1	1	1	1.0	0	8	16
1943	0.015348874	0.001534887	1.07274501	0.107274502	2.326659441	1240.774219	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	4	1	1	2	1.5	0	8	16
1944	0.007332874	0.000733287	0.285362167	0.028536217	2.397288656	1241.660156	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	2	1	4	1	0	1	1.0	0	8	8
1945	0.03894474	0.003894474	1.836211315	0.18362113	2.382192755	1241.624219	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	2	1	4	1	0	2	1.5	0	8	8
1946	0.000533205	5.33e-05	0.087070911	0.008707091	2.334634495	1240.799219	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	2	2	4	1	1	1	1.0	0	8	8
1947	0.008562821	0.000856282	1.201094464	0.120109444	2.31390729	1241.557813	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	2	2	4	1	1	2	1.5	0	8	8
1948	0.004823769	0.000482377	0.222542943	0.022254294	2.368333817	1240.717188	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	2	1	4	1	0	1	1.0	0	8	16
1949	0.031591403	0.00315914	1.51600726	0.151600727	2.371254349	1240.785938	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	2	1	4	1	0	2	1.5	0	8	16
1950	0.001338041	0.000133804	0.084246678	0.008424668	2.333894444	1240.533594	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	2	2	4	1	1	1	1.0	0	8	16
1951	0.015348874	0.001534887	1.07274501	0.107274502	2.265339231	1240.628906	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	2	2	4	1	1	2	1.5	0	8	16
1952	0.01089445	0.001089445	0.453577777	0.045357777	2.382671642	1240.747656	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	1	1	4	2	0	1	1.0	0	8	8
1953	0.047034564	0.004703456	2.075872362	0.207587238	2.319334459	1240.703906	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	1	1	4	2	0	2	1.5	0	8	8
1954	0.000278009	2.78e-05	0.088678933	0.008867893	2.290198708	1240.641406	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	4	2	1	1	1.0	0	8	8
1955	0.00825585	0.000825585	1.284150191	0.128415016	2.280206776	1240.852344	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	4	2	1	2	1.5	0	8	8
1956	0.008221407	0.000822141	0.366433837	0.036643384	2.366846514	1240.738281	0.0005	40	0.01	1	0.0005	0.0	0.08	const	0	0	0	3		1	5	1	1	4	2	0	1	1.0	0	8	16
1957	0.029226536	0.002922654	1.6304086	0.16304086	2.337304449	1240.703906	0.00025	40	0.02	1	0.0005	0.0	0.04	const	0	0	0	3		1	5	1	1	4	2	0	2	1.5	0	8	16
1958	0.002771273	0.000277127	0.116460946	0.011646095	2.293221283	1240.6875	0.0005	40	0.01	1	0.0005	0.8	0.08	harmonic	0	0	0	3	0.1	1	5	1	2	4	2	1	1	1.0	0	8	16
1959	0.012767993	0.001276799	1.113148995	0.1113149	2.383634329	1240.6625	0.00025	40	0.02	1	0.0005	0.8	0.04	harmonic	0	0	0	3	0.05	1	5	1	2	4	2	1	2	1.5	0	8	16

M ABLATION ON DECEPTIVE MEAN

We conducted an ablation study for the Deceptive Message game to analyze the sensitivity of GEMS to its core hyperparameters. Our investigation focused on key parameters governing the Amortized Best-Response (ABR) training, including the learning rate (`abr_lr`), update steps (`abr_steps`), and KL-divergence coefficient (`beta_kl`), as well as the generator’s latent dimension (`zdim`) and the Jacobian regularization penalty (`lambda_jac`).

Table 2: Ablation (Deceptive Mean)

	sender_last	receiver_last	sender_auc	receiver_auc	time_total_sec_mu	ram_peak_mb_mu	abr_lr	abr_steps	adv_ema	beta_kl	inner_eval	lambda_jac	meta_eval
--	-------------	---------------	------------	--------------	-------------------	----------------	--------	-----------	---------	---------	------------	------------	-----------

	sender_last	receiver_last	sender_auc	receiver_auc	time_total_sec_mu	ram_peak_mb_mu	abr Jr	abr_steps	adv_ema	beta_kl	inner_eval	lambda_jac	meta_eval_n	zdim
1972	0.028320313	0.71875	0.391349073	3.048147321	3.485627413	766.8632813	0.0005	100.0	0.1	0.01	4.0	0.0	32.0	16.0
1973	0.03125	0.703125	0.41984623	3.008728387	4.091207623	766.8574219	0.0005	100.0	0.05	0.01	2.0	0.0	64.0	8.0
1974	0.029296875	0.70703125	0.433813333	3.062614619	4.114554763	767.0214844	0.0005	100.0	0.05	0.01	2.0	0.0	64.0	16.0
1975	0.03125	0.705078125	0.41971549	3.019714162	4.117322803	766.9511719	0.0005	100.0	0.1	0.01	2.0	0.0	64.0	8.0
1976	0.032226563	0.704101563	0.440840818	3.053851022	4.145797491	766.9238281	0.0005	100.0	0.1	0.01	2.0	0.0	64.0	16.0
1977	0.03125	0.703125	0.41984623	3.008728387	4.09178257	766.828125	0.0005	100.0	0.05	0.01	4.0	0.0	64.0	8.0
1978	0.029296875	0.70703125	0.433813333	3.062614619	4.104922652	766.8183594	0.0005	100.0	0.05	0.01	4.0	0.0	64.0	16.0
1979	0.03125	0.705078125	0.41971549	3.019714162	4.11270225	766.8984375	0.0005	100.0	0.1	0.01	4.0	0.0	64.0	8.0
1980	0.032226563	0.704101563	0.440840818	3.053851022	4.102519035	766.8339844	0.0005	100.0	0.1	0.01	4.0	0.0	64.0	16.0
1981	0.03125	0.703125	0.41984623	3.008728387	4.09178257	766.828125	0.0005	100.0	0.05	0.01	4.0	0.0	64.0	8.0
1982	0.029296875	0.70703125	0.433813333	3.062614619	4.104922652	766.8183594	0.0005	100.0	0.05	0.01	4.0	0.0	64.0	16.0
1983	0.03125	0.705078125	0.41971549	3.019714162	4.11270225	766.8984375	0.0005	100.0	0.1	0.01	4.0	0.0	64.0	8.0
1984	0.032226563	0.704101563	0.440840818	3.053851022	4.102519035	766.8339844	0.0005	100.0	0.1	0.01	4.0	0.0	64.0	16.0
1985	0.002929688	0.783203125	0.081438271	3.804808408	4.592607737	766.9003906	0.001	200.0	0.05	0.01	2.0	0.0	32.0	8.0
1986	0.004882813	0.779296875	0.103168447	3.670696127	4.610334158	766.8359375	0.001	200.0	0.05	0.01	2.0	0.0	32.0	16.0
1987	0.006835938	0.78125	0.078946951	3.813918651	4.574136138	766.8300781	0.001	200.0	0.1	0.01	2.0	0.0	32.0	8.0
1988	0.004882813	0.778320313	0.094999513	3.680365159	4.598544359	766.796875	0.001	200.0	0.1	0.01	2.0	0.0	32.0	16.0
1989	0.002929688	0.783203125	0.081438271	3.804808408	4.564905047	766.8046875	0.001	200.0	0.05	0.01	4.0	0.0	32.0	8.0
1990	0.004882813	0.779296875	0.103168447	3.670696127	4.548526252	766.9375	0.001	200.0	0.05	0.01	4.0	0.0	32.0	16.0
1991	0.006835938	0.78125	0.078946951	3.813918651	4.592769027	767.0585938	0.001	200.0	0.1	0.01	4.0	0.0	32.0	8.0
1992	0.004882813	0.778320313	0.094999513	3.680365159	4.618440151	766.921875	0.001	200.0	0.1	0.01	4.0	0.0	32.0	16.0
1993	0.0	0.805664063	0.02203125	3.879985872	5.202624202	767.0332031	0.001	200.0	0.05	0.01	2.0	0.0	64.0	8.0
1994	0.002929688	0.791992188	0.075109614	3.907348489	5.187786818	766.7890625	0.001	200.0	0.05	0.01	2.0	0.0	64.0	16.0
1995	0.0	0.8046875	0.020729167	3.888508007	5.156398416	766.8984375	0.001	200.0	0.1	0.01	2.0	0.0	64.0	8.0
1996	0.002929688	0.791992188	0.075109614	3.907348489	5.170324206	766.9003906	0.001	200.0	0.1	0.01	2.0	0.0	64.0	16.0
1997	0.0	0.805664063	0.02203125	3.879985872	5.196808338	766.9003906	0.001	200.0	0.05	0.01	4.0	0.0	64.0	8.0
1998	0.002929688	0.791992188	0.075109614	3.907348489	5.195022225	766.8632813	0.001	200.0	0.05	0.01	4.0	0.0	64.0	16.0
1999	0.0	0.8046875	0.020729167	3.888508007	5.201624393	767.0410156	0.001	200.0	0.1	0.01	4.0	0.0	64.0	8.0
2000	0.002929688	0.791992188	0.075109614	3.907348489	5.167958856	766.84375	0.001	200.0	0.1	0.01	4.0	0.0	64.0	16.0
2001	0.013671875	0.759765625	0.233534138	3.491052739	4.574347615	766.8769531	0.0005	200.0	0.05	0.01	2.0	0.0	32.0	8.0
2002	0.008789063	0.760742188	0.237348755	3.381854229	4.594873428	767.046875	0.0005	200.0	0.05	0.01	2.0	0.0	32.0	16.0
2003	0.015625	0.759765625	0.2345107	3.502512135	4.615203023	766.9511719	0.0005	200.0	0.1	0.01	2.0	0.0	32.0	8.0
2004	0.01171875	0.755859375	0.240472116	3.358276201	4.573800206	766.9648438	0.0005	200.0	0.1	0.01	2.0	0.0	32.0	16.0
2005	0.013671875	0.759765625	0.233534138	3.491052739	4.576252937	766.8632813	0.0005	200.0	0.05	0.01	4.0	0.0	32.0	8.0
2006	0.008789063	0.760742188	0.237348755	3.381854229	4.60927546	767.0957031	0.0005	200.0	0.05	0.01	4.0	0.0	32.0	16.0
2007	0.015625	0.759765625	0.2345107	3.502512135	4.562783718	766.9199219	0.0005	200.0	0.1	0.01	4.0	0.0	32.0	8.0
2008	0.01171875	0.755859375	0.240472116	3.358276201	4.630032897	766.8847656	0.0005	200.0	0.1	0.01	4.0	0.0	32.0	16.0
2009	0.00390625	0.797851563	0.123591093	3.614750257	5.193202019	766.8613281	0.0005	200.0	0.05	0.01	2.0	0.0	64.0	8.0
2010	0.002929688	0.7890625	0.241208324	3.576210672	5.175317168	766.8652344	0.0005	200.0	0.05	0.01	2.0	0.0	64.0	16.0
2011	0.00390625	0.796875	0.125327204	3.619904337	5.203407407	766.9238281	0.0005	200.0	0.1	0.01	2.0	0.0	64.0	8.0
2012	0.002929688	0.787109375	0.252284713	3.56288221	5.175059795	766.9140625	0.0005	200.0	0.1	0.01	2.0	0.0	64.0	16.0
2013	0.00390625	0.797851563	0.123591093	3.614750257	5.206146717	766.9824219	0.0005	200.0	0.05	0.01	4.0	0.0	64.0	8.0
2014	0.002929688	0.7890625	0.241208324	3.576210672	5.198497653	766.8398438	0.0005	200.0	0.05	0.01	4.0	0.0	64.0	16.0
2015	0.00390625	0.796875	0.125327204	3.619904337	5.194039464	766.8378906	0.0005	200.0	0.1	0.01	4.0	0.0	64.0	8.0
2016	0.002929688	0.787109375	0.252284713	3.56288221	5.209975243	766.8867188	0.0005	200.0	0.1	0.01	4.0	0.0	64.0	16.0
2017	0.004882813	0.750976563	0.209971877	3.645173301	4.484597445	766.8574219	0.001	100.0	0.05	0.005	2.0	0.0	32.0	8.0
2018	0.001953125	0.7734375	0.171805644	3.671331845	3.481884599	766.9511719	0.001	100.0	0.05	0.005	2.0	0.0	32.0	16.0
2019	0.00390625	0.75390625	0.200136232	3.649067637	3.479956746	766.9042969	0.001	100.0	0.1	0.005	2.0	0.0	32.0	8.0
2020	0.001953125	0.774414063	0.15623379	3.687826008	3.488405466	766.7949219	0.001	100.0	0.1	0.005	2.0	0.0	32.0	16.0
2021	0.004882813	0.750976563	0.209971877	3.645173301	3.495933414	766.8652344	0.001	100.0	0.05	0.005	4.0	0.0	32.0	8.0
2022	0.001953125	0.7734375	0.171805644	3.671331845	3.493604779	766.9453125	0.001	100.0	0.05	0.005	4.0	0.0	32.0	16.0
2023	0.00390625	0.75390625	0.200136232	3.649067637	3.50344944	766.8417969	0.001	100.0	0.1	0.005	4.0	0.0	32.0	8.0
2024	0.001953125	0.774414063	0.15623379	3.687826008	3.479611635	766.8710938	0.001	100.0	0.1	0.005	4.0	0.0	32.0	16.0
2025	0.005859375	0.764648438	0.185098763	3.501440219	4.079826832	766.8730469	0.001	100.0	0.05	0.005	2.0	0.0	64.0	8.0
2026	0.004882813	0.765625	0.19841044	3.547859269	4.141643167	766.9648438	0.001	100.0	0.05	0.005	2.0	0.0	64.0	16.0
2027	0.005859375	0.765625	0.177693718	3.505019841	4.081937313	766.8808594	0.001	100.0	0.1	0.005	2.0	0.0	64.0	8.0
2028	0.0078125	0.760742188	0.192957412	3.55072115	4.103323817	766.8632813	0.001	100.0	0.1	0.005	2.0	0.0	64.0	16.0
2029	0.005859375	0.764648438	0.185098763	3.501440219	4.093967438	766.8730469	0.001	100.0	0.05	0.005	4.0	0.0	64.0	8.0
2030	0.004882813	0.765625	0.19841044	3.547859269	4.092059493	766.9082031	0.001	100.0	0.05	0.005	4.0	0.0	64.0	16.0
2031	0.005859375	0.765625	0.177693718	3.505019841	4.076184273	766.9199219	0.001	100.0	0.1	0.005	4.0	0.0	64.0	8.0
2032	0.0078125	0.760742188	0.192957412	3.55072115	4.161067963	767.0039063	0.001	100.0	0.1	0.005	4.0	0.0	64.0	16.0
2033	0.029296875	0.688476563	0.393779673	3.103399988	3.47068882	766.890625	0.0005	100.0	0.05	0.005	2.0	0.0	32.0	8.0
2034	0.02734375	0.712890625	0.411155931	3.030212674	3.492367029	766.9921875	0.0005	100.0	0.05	0.005	2.0	0.0	32.0	16.0
2035	0.03125	0.6875	0.397716837	3.114880598	3.490564466	766.9101563	0.0005	100.0	0.1	0.005	2.0	0.0	32.0	8.0
2036	0.028320313	0.719726563	0.391349073	3.054277964	3.492498398	766.8339844	0.0005	100.0	0.1	0.005	2.0	0.0	32.0	16.0
2037	0.029296875	0.688476563	0.393779673	3.103399988	3.50611043	766.8066406	0.0005	100.0	0.05	0.005	4.0	0.0	32.0	8.0
2038	0.02734375	0.712890625	0.411155931	3.030212674	3.485206127	766.7988281	0.0005	100.0	0.05	0.005	4.0	0.0	32.0	16.0
2039	0.03125	0.6875	0.397716837	3.114880598	3.504174113	766.8613281	0.0005	100.0	0.1	0.005	4.0	0.0	32.	

	sender_last	receiver_last	sender_auc	receiver_auc	time_total_sec_mu	ram_peak_mb_mu	abr Jr	abr_steps	adv_ema	beta_kl	inner_eval	lambda_jac	meta_eval_n	zdim
2030														
2031														
2032														
2033														
2034														
2035														
2036	0.004882813	0.783203125	0.071694878	3.816934612	4.575098038	766.9042969	0.001	200.0	0.1	0.005	4.0	0.0	32.0	8.0
2037	0.0	0.806640625	0.020295139	3.88961531	5.265115499	766.9238281	0.001	200.0	0.05	0.005	2.0	0.0	64.0	8.0
2038	0.0	0.793945313	0.07736926	3.900982364	5.19207108	766.8710938	0.001	200.0	0.05	0.005	2.0	0.0	64.0	16.0
2039	0.0	0.8046875	0.016822917	3.885627126	5.201019526	766.9609375	0.001	200.0	0.1	0.005	2.0	0.0	64.0	8.0
2040	0.0	0.793945313	0.065121528	3.921266165	5.193174481	766.8808594	0.001	200.0	0.1	0.005	2.0	0.0	64.0	16.0
2041	0.0	0.806640625	0.020295139	3.88961531	5.191734433	766.9140625	0.001	200.0	0.05	0.005	4.0	0.0	64.0	8.0
2042	0.0	0.793945313	0.07736926	3.900982364	5.205765486	766.9570313	0.001	200.0	0.05	0.005	4.0	0.0	64.0	16.0
2043	0.0	0.8046875	0.016822917	3.885627126	5.192848086	766.8886719	0.001	200.0	0.1	0.005	4.0	0.0	64.0	8.0
2044	0.0	0.793945313	0.065121528	3.921266165	5.187784672	767.0644531	0.001	200.0	0.1	0.005	4.0	0.0	64.0	16.0
2045	0.012695313	0.762695313	0.231309745	3.499253694	4.587309361	766.8378906	0.0005	200.0	0.05	0.005	2.0	0.0	32.0	8.0
2046	0.009765625	0.76171875	0.233089303	3.88961531	4.590754867	766.9121094	0.0005	200.0	0.05	0.005	2.0	0.0	32.0	16.0
2047	0.012695313	0.76171875	0.234781967	3.486716801	4.591512561	766.9394531	0.0005	200.0	0.1	0.005	2.0	0.0	32.0	8.0
2048	0.010742188	0.760742188	0.233577585	3.382513951	4.563710451	767.03125	0.0005	200.0	0.1	0.005	2.0	0.0	32.0	16.0
2049	0.012695313	0.762695313	0.231309745	3.499253694	4.560121536	766.859375	0.0005	200.0	0.05	0.005	4.0	0.0	32.0	8.0
2050	0.009765625	0.76171875	0.233089303	3.394772003	4.592621326	766.8652344	0.0005	200.0	0.05	0.005	4.0	0.0	32.0	16.0
2051	0.012695313	0.76171875	0.234781967	3.486716801	4.588880777	766.8007813	0.0005	200.0	0.1	0.005	4.0	0.0	32.0	8.0
2052	0.010742188	0.760742188	0.233577585	3.382513951	4.575504876	766.9824219	0.0005	200.0	0.1	0.005	4.0	0.0	32.0	16.0
2053	0.00390625	0.797851563	0.121854982	3.611278035	5.235417843	766.8730469	0.0005	200.0	0.05	0.005	2.0	0.0	64.0	8.0
2054	0.000976563	0.7890625	0.23594463	3.584682894	5.163270712	766.9960938	0.0005	200.0	0.05	0.005	2.0	0.0	64.0	16.0
2055	0.00390625	0.797851563	0.121854982	3.611278035	5.172132373	766.859375	0.0005	200.0	0.1	0.005	2.0	0.0	64.0	8.0
2056	0.000976563	0.7890625	0.23594463	3.584682894	5.199103951	766.921875	0.0005	200.0	0.05	0.005	4.0	0.0	64.0	16.0
2057	0.00390625	0.797851563	0.125327204	3.622892618	5.173853517	766.9179688	0.0005	200.0	0.1	0.005	4.0	0.0	64.0	8.0
2058	0.001953125	0.788085938	0.239114672	3.573882113	5.201363683	766.8847656	0.0005	200.0	0.1	0.005	2.0	0.0	64.0	16.0
2059	0.00390625	0.797851563	0.121854982	3.611278035	5.172132373	766.9941406	0.0005	200.0	0.05	0.005	4.0	0.0	64.0	8.0
2060	0.000976563	0.7890625	0.23594463	3.584682894	5.199103951	766.921875	0.0005	200.0	0.05	0.005	4.0	0.0	64.0	16.0
2061	0.00390625	0.797851563	0.125327204	3.622892618	5.173853517	766.9179688	0.0005	200.0	0.1	0.005	4.0	0.0	64.0	8.0
2062	0.001953125	0.788085938	0.239114672	3.573882113	5.177394509	766.8203125	0.0005	200.0	0.1	0.005	4.0	0.0	64.0	16.0
2063	0.008789063	0.73828125	0.234759823	3.582373689	3.485162735	766.953125	0.001	100.0	0.05	0.01	2.0	0.01	32.0	8.0
2064	0.009765625	0.76171875	0.214884584	3.61147658	3.503759861	766.8085938	0.001	100.0	0.05	0.01	2.0	0.01	32.0	16.0
2065	0.015625	0.73046875	0.244402282	3.574820543	3.476381302	766.8730469	0.001	100.0	0.1	0.01	2.0	0.01	32.0	8.0
2066	0.002929688	0.771484375	0.165628322	3.683584449	3.504372954	766.9472656	0.001	100.0	0.1	0.01	2.0	0.01	32.0	16.0
2067	0.008789063	0.73828125	0.234759823	3.582373689	3.492276192	766.9433594	0.001	100.0	0.05	0.01	4.0	0.01	32.0	8.0
2068	0.009765625	0.76171875	0.214884584	3.61147658	3.477153182	766.8515625	0.001	100.0	0.05	0.01	4.0	0.01	32.0	16.0
2069	0.015625	0.73046875	0.244402282	3.574820543	3.462974429	766.7929688	0.001	100.0	0.1	0.01	4.0	0.01	32.0	8.0
2070	0.002929688	0.771484375	0.165628322	3.683584449	3.494990706	766.9550781	0.001	100.0	0.1	0.01	4.0	0.01	32.0	16.0
2071	0.00390625	0.763671875	0.186622201	3.495951938	4.06504941	766.9042969	0.001	100.0	0.05	0.01	2.0	0.01	64.0	8.0
2072	0.005859375	0.762695313	0.190373352	3.543969893	4.111203074	766.9824219	0.001	100.0	0.05	0.01	2.0	0.01	64.0	16.0
2073	0.00390625	0.766601563	0.176413868	3.512755855	4.090454459	766.8203125	0.001	100.0	0.1	0.01	2.0	0.01	64.0	8.0
2074	0.0078125	0.752929688	0.21913159	3.464269195	4.129907131	766.9765625	0.001	100.0	0.1	0.01	2.0	0.01	64.0	16.0
2075	0.00390625	0.763671875	0.186622201	3.495951938	4.090277076	766.8945313	0.001	100.0	0.05	0.01	4.0	0.01	64.0	8.0
2076	0.005859375	0.762695313	0.190373352	3.543969893	4.139833689	766.9960938	0.001	100.0	0.05	0.01	4.0	0.01	64.0	16.0
2077	0.00390625	0.766601563	0.176413868	3.512755855	4.094886661	766.9941406	0.001	100.0	0.1	0.01	4.0	0.01	64.0	8.0
2078	0.0078125	0.752929688	0.21913159	3.464269195	4.105113029	766.828125	0.001	100.0	0.1	0.01	4.0	0.01	64.0	16.0
2079	0.01953125	0.711914063	0.406643105	3.163368808	3.509097457	766.8496094	0.0005	100.0	0.05	0.01	2.0	0.01	32.0	8.0
2080	0.043945313	0.674804688	0.488205295	2.925471717	3.498527765	766.9296875	0.0005	100.0	0.05	0.01	2.0	0.01	32.0	16.0
2081	0.020507813	0.708984375	0.409682407	3.161523083	3.485819697	766.8867188	0.0005	100.0	0.1	0.01	2.0	0.01	32.0	8.0
2082	0.041992188	0.684570313	0.466933594	2.957788983	3.488726735	766.9667969	0.0005	100.0	0.1	0.01	2.0	0.01	32.0	16.0
2083	0.01953125	0.711914063	0.406643105	3.163368808	3.499094486	766.8027344	0.0005	100.0	0.05	0.01	4.0	0.01	32.0	8.0
2084	0.043945313	0.674804688	0.488205295	2.925471717	3.48725915	767.0683594	0.0005	100.0	0.05	0.01	4.0	0.01	32.0	16.0
2085	0.020507813	0.708984375	0.409682407	3.161523083	3.473590493	767.0058594	0.0005	100.0	0.1	0.01	4.0	0.01	32.0	8.0
2086	0.041992188	0.684570313	0.466933594	2.957788983	3.511683583	766.9941406	0.0005	100.0	0.1	0.01	4.0	0.01	32.0	16.0
2087	0.03125	0.713867188	0.422857851	3.029866204	4.102545857	766.9316406	0.0005	100.0	0.05	0.01	2.0	0.01	64.0	8.0
2088	0.037109375	0.69921875	0.467807097	3.024323448	4.104898214	767.0546875	0.0005	100.0	0.05	0.01	2.0	0.01	64.0	16.0
2089	0.030273438	0.716796875	0.416727209	3.04256475	4.077229142	766.9355469	0.0005	100.0	0.1	0.01	2.0	0.01	64.0	8.0
2090	0.040039063	0.694335938	0.483740832	3.006518521	4.094441156	766.8378906	0.0005	100.0	0.1	0.01	2.0	0.01	64.0	16.0
2091	0.03125	0.713867188	0.422857851	3.029866204	4.107426643	766.9726563	0.0005	100.0	0.05	0.01	4.0	0.01	64.0	8.0
2092	0.037109375	0.69921875	0.467807097	3.024323448	4.088187575	767.1386719	0.0005	100.0	0.05	0.01	4.0	0.01	64.0	16.0
2093	0.030273438	0.716796875	0.416727209	3.04256475	4.133657098	766.84375	0.0005	100.0	0.1	0.01	4.0	0.01	64.0	8.0
2094	0.040039063	0.694335938	0.483740832	3.006518521	4.127341628	766.84375	0.0005	100.0	0.1	0.01	4.0	0.01	64.0	16.0
2095	0.0078125	0.779296875	0.086942319	3.802173239	4.544541717	766.9042969	0.001	200.0	0.05	0.01	2.0	0.01	32.0	8.0
2096	0.008789063	0.772460938	0.139713054	3.637912371	4.565919995	766.9628906	0.001	200.0	0.05	0.01	2.0	0.01	32.0	16.0
2097	0.002929688	0.783203125	0.073504243	3.863313226	4.57320559	766.9121094	0.001	200.0	0.1	0.01	2.0	0.01	32.0	8.0
2098	0.001953125	0.784179688	0.126850641	3.663778832	4.571519375	766.9609375	0.001	200.0	0.1	0.01	2.0	0.01	32.0	16.0
2099	0.0078125	0.779296875	0.086942319	3.802173239	4.599365115	766.90625	0.001	200.0	0.05	0.01	4.0	0.01	32.0	8.0
2100	0.008789063	0.772460938	0.139713054	3.637912371	4.587291718	766.8632813	0.001	200.0	0.05	0.01	4.0	0.01	32.0	16.0
2101	0.002929688	0												

2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145

sender_last	receiver_last	sender_auc	receiver_auc	time_total_sec_mu	ram_peak_mb_mu	abr Jr	abr_steps	adv_ema	beta_kl	inner_eval	lambda_jac	meta_eval_n	zdim
0.006835938	0.767578125	0.259089737	3.39466137	4.57443428	766.9101563	0.0005	200.0	0.05	0.01	4.0	0.01	32.0	16.0
0.008789063	0.774414063	0.240667916	3.534573811	4.563651443	766.8945313	0.0005	200.0	0.1	0.01	4.0	0.01	32.0	8.0
0.006835938	0.764648438	0.259012144	3.382884026	4.570395947	766.8925781	0.0005	200.0	0.1	0.01	4.0	0.01	32.0	16.0
0.004882813	0.791992188	0.133063217	3.590356213	5.188795686	766.8652344	0.0005	200.0	0.05	0.01	2.0	0.01	64.0	8.0
0.005859375	0.774414063	0.252983542	3.493388295	5.219975233	766.90625	0.0005	200.0	0.05	0.01	2.0	0.01	64.0	16.0
0.00390625	0.793945313	0.131299426	3.608049399	5.20837307	766.8339844	0.0005	200.0	0.1	0.01	2.0	0.01	64.0	8.0
0.004882813	0.7734375	0.253770771	3.485601261	5.190756202	766.8964844	0.0005	200.0	0.1	0.01	2.0	0.01	64.0	16.0
0.004882813	0.791992188	0.133063217	3.590356213	5.159092546	766.9257813	0.0005	200.0	0.05	0.01	4.0	0.01	64.0	8.0
0.005859375	0.774414063	0.252983542	3.493388295	5.208868861	766.9707031	0.0005	200.0	0.05	0.01	4.0	0.01	64.0	16.0
0.00390625	0.793945313	0.131299426	3.608049399	5.236365676	766.9394531	0.0005	200.0	0.1	0.01	4.0	0.01	64.0	8.0
0.004882813	0.7734375	0.253770771	3.485601261	5.219882965	766.9980469	0.0005	200.0	0.1	0.01	4.0	0.01	64.0	16.0
0.005859375	0.744140625	0.224260116	3.612463241	3.47030735	766.8144531	0.001	100.0	0.05	0.005	2.0	0.01	32.0	8.0
0.010742188	0.76171875	0.215372865	3.617039222	3.490041375	766.8730469	0.001	100.0	0.05	0.005	2.0	0.01	32.0	16.0
0.017578125	0.727539063	0.245378844	3.584327921	3.494503975	766.8828125	0.001	100.0	0.1	0.005	2.0	0.01	32.0	8.0
0.00390625	0.770507813	0.166116603	3.689371678	3.497444153	766.9824219	0.001	100.0	0.1	0.005	2.0	0.01	32.0	16.0
0.005859375	0.744140625	0.224260116	3.612463241	3.486721992	766.7871094	0.001	100.0	0.05	0.005	4.0	0.01	32.0	8.0
0.010742188	0.76171875	0.215372865	3.617039222	3.518079996	766.8613281	0.001	100.0	0.05	0.005	4.0	0.01	32.0	16.0
0.017578125	0.727539063	0.245378844	3.584327921	3.484863997	766.8378906	0.001	100.0	0.1	0.005	4.0	0.01	32.0	8.0
