

# Evolutionary Policy Optimization

Jianren Wang<sup>\*1</sup> Yifan Su<sup>\*1</sup> Abhinav Gupta<sup>1</sup> Deepak Pathak<sup>1</sup>

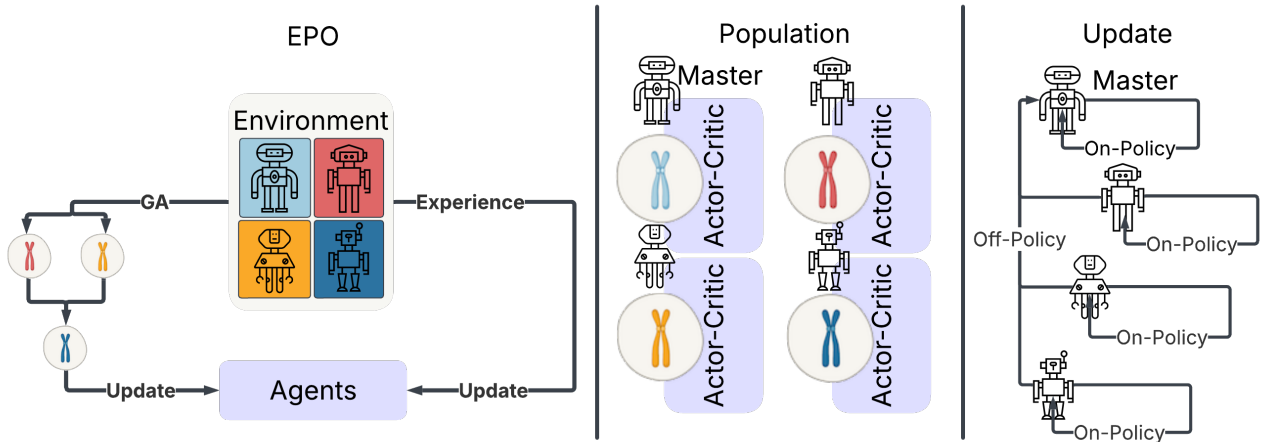


Figure 1. We introduce Evolutionary Policy Optimization (EPO) (Left), a novel policy gradient algorithm that integrates a genetic algorithm (GA) with policy gradients. During training, a population of agents, each represented by a latent embedding (gene) and sharing a common actor-critic network, interacts with the environment (Middle). Among them, one agent is designated as the master agent, which is trained using amortized experience aggregated from all agents in the population to enhance learning efficiency and stability (Right).

## Abstract

Despite its extreme sample inefficiency, on-policy reinforcement learning has become a fundamental tool in real-world applications. With recent advances in GPU-driven simulation, the ability to collect vast amounts of data for RL training has scaled exponentially. However, studies show that current on-policy methods, such as PPO, fail to fully leverage the benefits of parallelized environments, leading to performance saturation beyond a certain scale. In contrast, Evolutionary Algorithms (EAs) excel at increasing diversity through randomization, making them a natural complement to RL. However, existing EvoRL methods have struggled to gain widespread adoption due to their extreme sample inefficiency. To address these challenges, we introduce Evolutionary Policy Optimization (EPO), a novel policy gradient algorithm that combines the strengths of

EA and policy gradients. We show that EPO significantly improves performance across diverse and challenging environments, demonstrating superior scalability with parallelized simulations. For visualizations of the learned policies, please visit: <https://sites.google.com/view/epo-rl>.

## 1. Introduction

Reinforcement Learning (RL) has emerged as a powerful framework for training autonomous decision-making agents in various domains, including games [36, 2, 42, 27], robotics [17, 1, 20, 15], and large language models (LLMs) [29, 12, 41]. RL, whether off-policy or on-policy, follows a trial-and-error approach and is inherently sample inefficient. On-policy RL methods, in particular, suffer more from inefficiency compared to off-policy or model-based alternatives since they discard past experience.

Despite their sample inefficiency, on-policy model-free RL methods, also known as policy gradients [34, 27], are widely adopted in real-world applications [36, 2, 17, 1, 29]. Their popularity stems from their ability to achieve higher asymptotic rewards and stable performance, making them the preferred choice in data-rich environments such as simulations

<sup>\*</sup>Equal contribution, order decided by coin flip. <sup>1</sup>Robotics Institute, Carnegie Mellon University, PA15213, USA. Correspondence to: Jianren Wang <jianrenw@andrew.cmu.edu>, Yifan Su <yifansu2003@gmail.com>.

and game engines [23], where large training batches can be easily generated. However, recent studies [37, 30] reveal a critical limitation: on-policy RL methods do not scale well with increasing batch sizes. Beyond a certain point, performance plateaus despite additional data collection, indicating that simply generating more samples across GPUs does not necessarily improve learning.

This limitation arises from the fundamental nature of on-policy RL: since data is collected from the **current** policy, increasing the number of parallel environments does not guarantee more diverse experiences. Instead, data distributions tend to converge, reducing the benefits of additional sampling. Evolutionary Algorithms (EAs), in contrast, excel at increasing diversity through randomization, making them a natural complement to RL. Evolutionary Reinforcement Learning (EvoRL), which integrates EAs with RL, is not a new concept [28] and has shown promise in prior work [32]. By maintaining a population of diverse policies and refining them through evolutionary operations like crossover and mutation, EvoRL enhances policy search and exploration.

However, EvoRL has not gained mainstream adoption due to its extreme sample inefficiency. In its simplest form, EvoRL resembles random search [24], as it lacks direct reward optimization like policy gradients. Typical EvoRL methods optimize reward indirectly through gradient-free search, known as Evolutionary Strategies (ES) [18, 3, 43], but training large deep models with ES is computationally prohibitive due to excessive memory and compute requirements. Even hybrid approaches that update a master agent via gradient ascent [31] struggle with efficiently maintaining and updating a large population of agents within GPU constraints, leading to low-quality data generation and poor sample efficiency. A notable hybrid approach, Population-Based Training (PBT) [30], aims to combine policy gradients with EvoRL by evolving a population of PPO policies. However, PBT segregates each policy in the population, preventing effective amortization of experience across agents and ultimately limiting its performance.

In this paper, we introduce Evolutionary Policy Optimization (EPO), a novel policy gradient algorithm that combines the strengths of EA and policy gradients: 1) EA to increase the diversity of experience in massively parallelized training 2) Amortized experience from all policies in a population to directly train a master policy. To achieve this, we adopt the split and aggregate policy gradients (SAPG) formulation [37], which performs off-policy updates via importance sampling to aggregate experience from individual follower policies into the on-policy data of the master policy.

EPO maintains a population of policies alongside a central master policy, all sharing the same network weights (analogous to gene-expression rules). This prevents parameter bloat as the population scales, enabling efficient training

and storage of large networks. Each agent is conditioned on a unique latent variable (“gene”), ensuring behavioral diversity while maintaining coherence across policies. Periodically, Darwinian natural selection [7] is applied: low-performing agents are eliminated, while top-performing elite agents are preserved. These elites undergo crossover and mutation, injecting controlled diversity while avoiding excessive divergence. All policies are updated synchronously via policy gradients, using shared reward signals to ensure efficient evolution, allowing EPO to scale effectively with larger neural networks. In each iteration, EPO aggregates experience from all policies into the master agent via SAPG updates [37], enabling the master policy to learn from diverse, high-quality data. This allows EPO to scale efficiently with increasing batch sizes, fully leveraging the potential of massively parallel simulations.

We evaluate EPO on several challenging RL benchmarks and demonstrate that it significantly outperforms prior state-of-the-art RL methods, including PBT [30] and SAPG [37]. Furthermore, we provide scaling laws, showing that EPO effectively scales with increased computational resources, making it a promising approach for large-scale RL training.

## 2. Related Work

**On-Policy Reinforcement Learning** On-policy approaches [39] update the policy using data collected by the current policy, ensuring that the training data remains closely aligned with the policy’s behavior. This strategy often leads to higher asymptotic rewards and stable performance, but at the cost of sample inefficiency, as past experiences from older policies are discarded. Despite this drawback, on-policy methods are the preferred choice in settings where data is abundant, such as simulations and game engines, where large training batches can be efficiently generated. In practice, on-policy RL has proven highly effective for training autonomous decision-making agents in games [36, 2, 42, 27], robotics [17, 1, 20, 15], and more recently, in developing reasoning capabilities in large language models (LLMs) [29, 12, 41]. Two widely used algorithms in this category are Trust Region Policy Optimization (TRPO) [33] and Proximal Policy Optimization (PPO) [34], both designed to stabilize learning by constraining the magnitude of policy updates.

**Evolutionary Reinforcement Learning** Another related line of work is Evolutionary Reinforcement Learning (EvoRL) [28], which integrates population-based search with policy optimization to explore a broad range of candidate solutions. Historically, evolutionary computation (EC) has been employed to optimize neural network weights [40, 4, 5], architectures [11], and hyperparameters [16, 38]. For instance, OpenAI ES [6] applied a stream-

lined evolution strategy to Atari games. More recently, researchers have integrated gradient-based updates alongside evolutionary operations (e.g., mutation and crossover), improving scalability and leveraging the representational power of deep neural networks [31, 19]. For example, [18] extends Deep Deterministic Policy Gradient (DDPG) by introducing evolutionary operations on a population of policies, combining the exploratory benefits of evolution with gradient-based optimization for more robust learning. However, these methods still struggle with efficiently maintaining and updating large agent populations within GPU constraints, leading to low-quality data generation and poor sample efficiency. In this work, we propose to combine the strengths of evolutionary search with on-policy updates, enabling higher asymptotic rewards while maintaining high sample efficiency.

**Off-Policy Reinforcement Learning** Unlike on-policy algorithms, off-policy methods allow for policy updates using data collected from different or older policies, enabling more efficient reuse of past experiences. Classic approaches such as Q-learning and its deep variant Deep Q-Network (DQN) [26] rely on a replay buffer to store state-action transitions for more stable learning. More recent continuous control algorithms, such as Deep Deterministic Policy Gradient (DDPG) [22, 35], Twin Delayed DDPG (TD3) [10], and Soft Actor-Critic (SAC) [13], extend these ideas within an actor-critic framework, further improving stability and performance. While off-policy methods are typically more sample-efficient than on-policy methods, they often suffer from distribution shift and overestimation of action values, which can hinder stable learning. To effectively amortize experience across all policies in a population while ensuring stable updates, we adopt the Split and Aggregate Policy Gradient (SAPG) formulation [37], which performs off-policy updates via importance sampling, enabling the aggregation of experience from individual follower policies into the on-policy data of a master policy.

### 3. Preliminaries

In this paper, we propose Evolutionary Policy Optimization (EPO), a novel reinforcement learning algorithm built on the well-established on-policy RL method Proximal Policy Optimization (PPO) and Genetic Algorithm (GA). Below, we provide a brief overview of these foundational methods.

**Reinforcement Learning** A standard reinforcement learning setting is formalized as a Markov Decision Process (MDP) and consists of an agent interacting with an environment  $E$  over a number of discrete time steps. At each time step  $t$ , the agent receives a state  $s_t$  and maps it to an action  $a_t$  using its policy  $\pi_\theta$ . The agent receives a scalar reward  $r_t$  and moves to the next state  $s_{t+1}$ . The process

continues until the agent reaches a terminal state marking the end of an episode. The return  $R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$  is the total accumulated return from time step  $t$  with discount factor  $\gamma \in (0, 1]$ . The goal of the agent is to maximize the expected return.

**PPO** Proximal Policy Optimization is widely used in reinforcement learning applications. It is an actor-critic based policy gradient algorithm, where the key idea underlying policy gradients is to increase the probabilities of actions that lead to higher return and decrease the probabilities of actions that lead to lower return, iteratively refining the policy to achieve optimal performance.

Let  $\pi_\theta$  denote a policy with parameters  $\theta$ , and  $J(\pi_\theta)$  denote the expected finite-horizon undiscounted return of the policy.  $J(\pi_\theta)$  is updated by:

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\substack{s \sim d^{\pi_\theta} \\ a \sim \pi_\theta(\cdot|s)}} [\nabla_\theta \log \pi_\theta(a|s) A^{\pi_\theta}(s, a)], \quad (1)$$

where  $A^{\pi_\theta}(s, a)$  is an advantage function that estimates the contribution of the transition to the gradient. A common choice is  $A^{\pi_\theta}(s, a) = Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s)$ , where  $Q^{\pi_\theta}(s, a)$  and  $V^{\pi_\theta}(s)$  are the estimated action-value and value functions, respectively. This form of update is termed an actor-critic update. Since we want the gradient of the error with respect to the current policy, only data from the current policy (on-policy) can be utilized. But these updates can be unstable because gradient estimates are high variance and the loss landscape is complex. An update step that is too large can degrade policy performance. Proximal Policy Optimization (PPO) modifies Eq. 1 to restrict updates to remain within an approximate *trust region* where improvement is guaranteed:

$$L_{on}(\pi_\theta) = \mathbb{E}_{\pi_{old}} [\min((r_t(\pi_\theta), \text{clip}(r_t(\pi_\theta), 1 - \epsilon, 1 + \epsilon)) A_t^{\pi_{old}})], \quad (2)$$

where  $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{old}(a_t|s_t)}$ ,  $\epsilon$  is a clipping hyperparameter, and  $\pi_{old}$  is the policy collecting the on-policy data. The clipping operation ensures that the updated  $\pi_\theta$  stays close to  $\pi_{old}$ . Empirically, given large numbers of samples, PPO achieves high performance and is stable and robust to hyperparameters.

**Genetic Algorithm (GA)** is one of the most well-known Evolutionary Algorithms (EAs). It relies on three key operators: selection, mutation, and crossover. A core concept in genetic algorithms is the population, where each individual represents a potential solution evaluated using a fitness function. The classic genetic algorithm emphasizes crossover as

the primary mechanism for exploration [14] and is usually applied to deal with the problem of hyper-parameter tuning in reinforcement learning.

## 4. Approach

The core idea behind Evolutionary Policy Optimization (EPO) is to integrate the population-based exploration of Genetic Algorithms (GAs) with the powerful policy gradient. While our framework specifically combines GA with Proximal Policy Optimization (PPO), it can be extended to any on-policy RL algorithm utilizing an actor-critic architecture.

EPO (As shown in Figure 1) begins by initializing a population of agents, where each agent is represented by a unique latent embedding (gene), while all agents share a common actor-critic network (gene-expression rules). Among these agents, one is designated as the “master” agent. Each agent then interacts with the environment for a full episode, and its fitness is determined by the cumulative reward obtained during that episode. A selection operator then determines which agents survive to the next generation, with survival probability proportional to their relative fitness scores. The embeddings of the surviving agents undergo mutation and crossover to form the next generation, while a small subset of top-performing agents (“elites”) is preserved unchanged to maintain stability.

To further enhance learning efficiency, we adopted a hybrid policy update strategy [37]. The master agent updates its parameters using both (i) on-policy data collected in the most recent episode and (ii) off-policy data sampled from the trajectories of the entire population. Meanwhile, each non-master agent is updated using only its own on-policy data. This hybrid learning process ensures both efficient exploration and sample reuse, making EPO well-suited for large-scale reinforcement learning. In the following sections, we describe each component of EPO in detail.

### 4.1. Genetic Algorithm for Agent Selection

The design philosophy of our algorithm is twofold: (1) to generate diverse yet bounded data for training the master agent and (2) to ensure that non-master agents are updated efficiently, enabling them to produce useful yet diverse data. Drawing inspiration from biological evolution, we propose initializing a population of agents ( $\{\pi_k\}_{k=1}^K$ ), where each agent is associated with a unique latent embedding ( $\{\phi_k\}_{k=1}^K$ ). All agents share a common actor-critic network, parameterized by  $\theta$  and  $\psi$ , respectively. The actor is conditioned on both the current state and its latent embedding to generate actions, while the value function is similarly conditioned to estimate value functions. Prior works have also leveraged latent-conditioned policies to encourage diverse behaviors [8, 9]. Among these agents, one is designated as

the master agent—for simplicity, we define the master agent as  $\pi_1(\theta, \psi, \phi_1)$ , though the general framework allows for alternative assignments.

To update non-master agents, we employ a genetic algorithm (GA) that modifies only the latent embeddings while keeping network parameters fixed. Each agent interacts with the environment for a full episode, and its fitness is determined by the cumulative reward obtained during that episode. A **selection** operator determines which agents survive to the next generation, with survival probability proportional to their relative fitness scores (*i.e.*, cumulative reward). Specifically, we retain the top  $x$  ranked non-master agents ( $x \in [2, K]$ ) as elites, forming a set  $X$ , which is preserved unchanged and carried over to the next generation set  $Y$ . We then randomly select an agent from both  $X$  and  $Y$  to undergo crossover and mutation. While multiple crossover methods exist (e.g., k-point crossover, uniform crossover), we use a simple averaging method for efficiency. Given two selected latent embeddings,  $\phi_i \in X$  and  $\phi_j \in Y$ , the **crossover** result is computed as:  $\phi' = 1/2 \times (\phi_i + \phi_j)$ . Following crossover, a **mutation** step is applied by adding Gaussian noise:  $\phi'' = \phi' + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma)$ . The mutated latent embedding  $\phi''$  is then added to  $Y$  for the next generation.

It is important to note that for more complex tasks, such as dexterous manipulation, policies require longer training before they exhibit meaningful differences in performance. In such cases (*e.g.*, when all rewards are close to zero), performing genetic updates too early may be ineffective. To address this, we execute selection, crossover, and mutation only when the agents’ performance differences become sufficiently significant. Specifically, the genetic algorithm is applied only when the difference between fitness scores exceeds a certain proportion of the median fitness score, ensuring that evolutionary updates occur only when meaningful behavioral differentiation has emerged.

$$\max\{f_k\}_{k=2}^K - \min\{f_k\}_{k=2}^K > \gamma \cdot \text{median}\{f_k\}_{k=2}^K \quad (3)$$

All parameters ( $\theta, \psi, \{\phi_k\}_{k=1}^K$ ) are updated using the following hybrid policy optimization process.

### 4.2. Hybrid-Policy Optimization

Different agents generate diverse data through varied behaviors, which can be beneficial if the master agent can effectively learn from all available data. However, since data distributions collected by different agents may vary significantly, incorporating off-policy data requires careful correction. Following [25, 37], we apply importance sampling to reweight updates from different policies. Specifically, we sample  $|S_1|$  transitions from  $\cup_{k=2}^K S_k$  collected by agents  $\{\pi_k\}_{k=2}^K$  to get  $S'_1$ . To update the master agent  $\pi_1$  using  $S'_1$ ,



we define the off-policy loss as:

$$L_{\text{off}}(\pi_1, S'_1) = \frac{1}{|S'_1|} \mathbb{E}_{(s,a) \sim S'_1} [\min(r_{\pi_1}(s, a), \text{clip}(r_{\pi_1}(s, a), \mu(1 - \epsilon), \mu(1 + \epsilon))) A^{\pi_1, \text{old}}(s, a)] \quad (4)$$

where  $r_{\pi_1}(s, a)$  is the importance sampling ratio, given by  $r_{\pi_1}(s, a) = \frac{\pi_1(s, a)}{\pi_k(s, a)}$  and  $\mu$  is an off-policy correction term, given by  $\mu = \frac{\pi_{1, \text{old}}(s, a)}{\pi_k(s, a)}$ . The final policy loss combines this off-policy loss with the on-policy loss given by Equation 2:

$$L(\pi_1) = L_{\text{on}}(\pi_1) + \lambda \cdot L_{\text{off}}(\pi_1, S'_1) \quad (5)$$

**Critic Update with Hybrid Data** The target value for the critic is computed using  $n$ -step returns (where  $n = 3$ ) for on-policy data:

$$\hat{V}_{\text{on}, \pi_k}^{\text{target}}(s_t) = \sum_{m=t}^{t+2} \gamma^{m-t} r_m + \gamma^3 V_{\pi_k, \text{old}}(s_{t+3}) \quad (6)$$

The critic loss for on-policy data is then:

$$L_{\text{on}}^{\text{critic}}(\pi_k) = \mathbb{E}_{(s,a) \sim \pi_k} \left[ \left( V_{\pi_k}(s) - \hat{V}_{\text{on}, \pi_k}^{\text{target}}(s) \right)^2 \right] \quad (7)$$

For off-policy data, however, multi-step returns are not directly applicable. Instead, we approximate the 1-step return:

$$\hat{V}_{\text{off}, \pi_1}^{\text{target}}(s'_t) = r_t + \gamma V_{\pi_1, \text{old}}(s'_{t+1}) \quad (8)$$

The critic loss for off-policy data is then:

$$L_{\text{off}}^{\text{critic}}(\pi_1, S'_1) = \frac{1}{|S'_1|} \mathbb{E}_{(s,a) \sim S'_1} \left[ \left( V_{\pi_1}(s) - \hat{V}_{\text{off}, \pi_1}^{\text{target}}(s) \right)^2 \right] \quad (9)$$

The final critic loss is a weighted combination of the on-policy and off-policy losses:

$$L^{\text{critic}}(\pi_1) = L_{\text{on}}^{\text{critic}}(\pi_1) + \lambda \cdot L_{\text{off}}^{\text{critic}}(\pi_1, S'_1) \quad (10)$$

**Updates for Non-Master Agents** All agents, except for the master agent  $\pi_1$ , are updated exclusively using on-policy losses, as defined in Equations 2 and 6. This ensures that while the master agent leverages additional experience from off-policy data, the remaining agents maintain stability by learning from their own trajectories, thereby preserving diversity within the population.

---

**Algorithm 1** Evolutionary Policy Optimization
 

---

```

1: Initialize  $K$  agents  $\{\pi_i\}_{i=1}^K \leftarrow (\{\phi_i\}_{i=1}^K, \theta, \psi)$ 
2: Initialize  $K$  empty cyclic replay buffers  $\{S_i\}_{i=1}^K$ 
3: Initialize  $N$  environments  $\{E\}_{i=1}^N$ 
4: for  $k = 1, \dots, K$  do
5:    $E^k = \{E\}_{i=\frac{(k-1)N}{K}}^{\frac{kN}{K}} \triangleright$  Assign environments
6: end for
7: for iteration=1,2,... do
8:   for  $k = 2, \dots, K$  do
9:      $f_k \leftarrow \text{Evaluate}(E^k, \theta, \psi, \phi_k)$ 
10:   end for
11:   if Eq 3 then
12:     Rank the population based on fitness scores  $f_k$ 
13:     Select top  $x$  as elites to form Set  $X$ , append to  $Y$ 
14:     while  $|Y| < K - 1$  do
15:       Crossover between a randomly sampled  $\phi_i \in X$ 
16:       and  $\phi_j \in Y$  to generate  $\phi'$ 
17:       Mutate( $\phi'$ ) and append it to  $Y$ 
18:     end while
19:   for  $k = 1, \dots, K$  do
20:      $R_k \leftarrow \text{CollectData}(E^k, \theta, \psi, \phi_k)$ 
21:   end for
22:    $L \leftarrow 0$ 
23:   Sample  $|S_1|$  transitions from  $\cup_{k=2}^K S_k$  to get  $S'_1$ 
24:    $L \leftarrow L + \text{OffPolicyLoss}(S'_1)$ 
25:   for  $k = 1, \dots, K$  do
26:      $L \leftarrow L + \text{OnPolicyLoss}(S_k)$ 
27:   end for
28:   Update  $\theta, \psi, \{\phi_k\}_{k=1}^K$ 
29: end for
30: Return  $\theta, \psi, \phi_1$ 
    
```

---

## 5. Experiments

We evaluate our method on four manipulation tasks, comparing it against state-of-the-art (SOTA) approaches. Additionally, we conduct ablation studies to assess the contribution of each component. To demonstrate scalability, we analyze our method’s performance with increasing computational resources. Finally, we investigate the impact of different crossover strategies. Each aspect is discussed in detail.

### 5.1. Experimental Setup

**Tasks** We evaluate our algorithm on four challenging dexterous manipulation tasks from the Isaac Gym Allegro-Kuka environments [23], including single-arm regrasping, single-arm reorientation, two-arm regrasping, and two-arm reorientation. These environments allow for diverse task completion strategies, such as throwing the object within a large exploration space, making them well-suited for evaluating the effectiveness of reinforcement learning algorithms. A

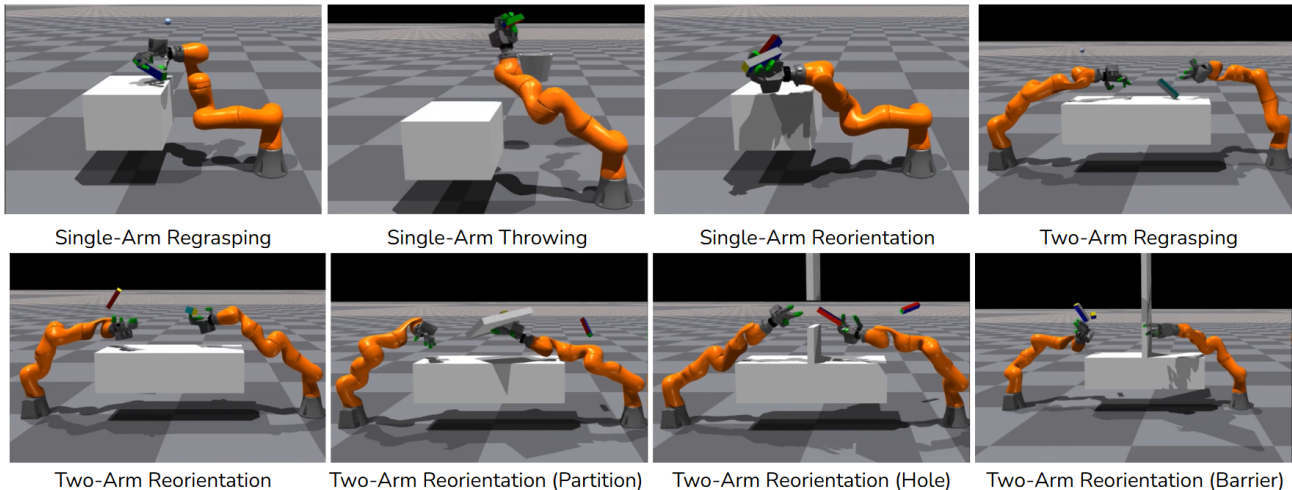


Figure 2. We evaluate our algorithm on four challenging dexterous manipulation tasks from the Isaac Gym Allegro-Kuka environments, including single-arm regrasping, single-arm reorientation, two-arm regrasping, and two-arm reorientation.

visual illustration of the tasks is provided in Figure 2.

**Evaluation** Following previous works [30, 37] we use the number of successes in a single episode as a performance metric on these tasks.

- **Single-Arm Regrasping:** The agent must lift a cube from the table and hold it near a randomly generated target point (shown as a white sphere in Figure 2) within 30 steps. A success is defined as the center distance between the cube and the target point being within a predefined threshold.
- **Single-Arm Throwing:** The agent need pick up the cube and throw it to a bucket at a random position. A success is defined as if the cube is thrown into the bucket.
- **Single-Arm Reorientation:** The agent must pick up an object and reorient it to a specified 6-DoF pose including both position and orientation (shown as a colored cube in Figure 2). A success is defined as both the positional and orientation differences between the object and the target pose being within a threshold.
- **Two-Arm Regrasping:** The goal is the same as in one-arm regrasping. However, the generated goal pose is far from the initial pose, requiring two Kuka arms to collaborate to transfer (*e.g.* throwing) the object to reach target poses in different regions.
- **Two-Arm Reorientation:** The goal is the same as in one-arm reorientation. However, the target pose is placed far from the initial pose, necessitating coordinated actions between the two Kuka arms.

- **Two-Arm Reorientation (Partition):** Same goal as two-arm reorientation, but there is a movable partition on the table. So, the agent should first learn to remove the partition and try to put the cube in the right position.
- **Two-Arm Reorientation (Hole):** Same goal as two-arm reorientation, but the barrier on the table has a hole in the middle. The agent should learn to throw the cube through the hole to successfully align the cube with its target.
- **Two-Arm Reorientation (Barrier):** Same goal as two-arm reorientation, but this time the barrier on the table cannot be passed through or removed. So one allegro hand must pass the cube from the side of the barrier to another hand.

**Baselines** We evaluate our approach against state-of-the-art RL methods specifically designed for the large-scale setting. Our comparisons include off-policy [21] and hybrid-policy [37] variants, population-based EvoRL [30], as well as the standard PPO [34].

- **Parallel Q-Learning [21]:** A parallelized version of DDPG [22], incorporating mixed exploration by varying exploration noise across environments to enhance exploration efficiency. This baseline is used to evaluate the scalability of off-policy algorithms as the number of samples increases. We use the official code for training without any algorithm-level modifications.
- **SAPG [37]:** A hybrid-policy RL algorithm designed for large-scale environments by splitting training data into chunks and recombining them via importance sampling. This baseline is used to assess the scalability

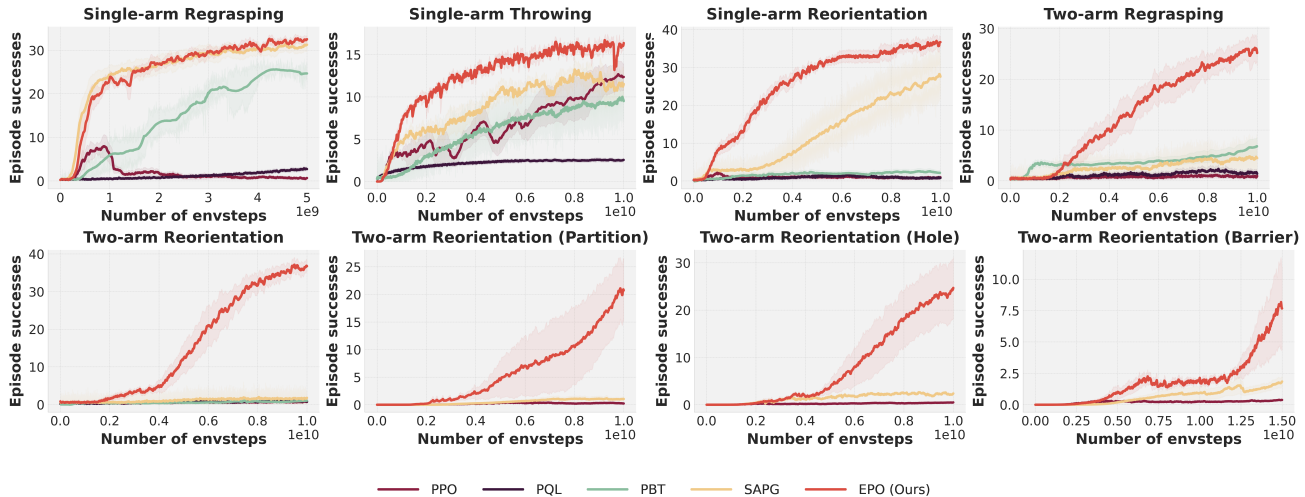


Figure 3. Performance curves of EPO with respect to PPO, PBT, PQL and SAPG baselines. In the experiment, EPO has 64 agents, and SAPG and PBT use the same number of agents as specified in their respective papers, namely 6 agents and 8 agents.

of on-policy algorithms as batch size increases. We use the official code for training without any algorithm-level modifications.

- **DexPBT [30]**: A population-based training (PBT) framework built on PPO. Similarly, the  $N$  environments are divided into  $K$  groups, each containing  $N/K$  environments, where  $K$  separate agents are trained with different hyperparameters. At regular intervals, the worst-performing policies are replaced. This baseline is used to evaluate the scalability of EvoRL. We use the official code for training without any algorithm-level modifications.
- **PPO [34]**: One of the most widely used reinforcement learning algorithms, serving as a strong baseline. We use this to demonstrate the scalability of standard RL.

For a fair comparison, all algorithms are trained using the same number of environments ( $N = 24,576$ ) and the same network architecture. To maintain consistency with the original papers, PBT is trained with 8 agents, and SAPG is trained with 6 agents, while EPO (ours) is trained with 64 agents. Additional ablation studies on the impact of the number of agents can be found in Section 5.3.

For simpler tasks such as regrasping, all algorithms are trained with  $5 \times 10^9$  transitions, while for more challenging tasks, all algorithms are trained with  $1 \times 10^{10}$  or  $1.5 \times 10^{10}$  transitions. Each experiment is run with 5 different random seeds, and we report the mean and standard error in the plots.

## 5.2. Results and Analysis

As shown in Figure 3 and Table 1, our algorithm consistently and significantly outperforms all baselines. For more visualization, please refer to <https://sites.google.com/view/epo-rl>.

For simpler tasks such as Kuka regrasping, SAPG and PBT perform comparably to our method, as the exploration space is relatively small. However, for more challenging tasks like single-arm reorientation, SAPG demonstrates some learning capability, whereas all other baselines fail. This highlights the advantage of incorporating genetic algorithms (GA), allowing our method to leverage large, high-quality populations to generate diverse and useful training data. Additionally, combining GA with the hybrid training strategy enables higher asymptotic performance compared to baseline methods.

For the two-arm tasks, where two Kuka robots must collaborate to complete the task, all baseline methods fail to learn any meaningful behavior even after  $1 \times 10^{10}$  steps, whereas EPO achieves substantial learning progress. In the original SAPG paper [37], the policy was trained for  $5 \times 10^{10}$  steps and converged at 28.58 episode successes, while EPO reaches 35.8 episode successes with only  $1 \times 10^{10}$  steps. This highlights the superior sample efficiency of EPO, significantly outperforming baseline methods. For the last three tasks, PPO and PQL are unable to learn any useful behaviors, so we omit their results in the figure and the result table.

Notably, EPO exhibits lower variance across different random seeds compared to baseline methods, with an average variance of 20.24% of the mean rewards, while the best-performing baseline reaches 37.50%. This indicates that

Task	SAPG[37]	PPO [34]	PBT [30]	PQL [21]	EPO (Ours)
Single Arm Regrasping	31.6 ± 2.07	0.1 ± 0.11	24.8 ± 1.17	3.1 ± 0.31	<b>32.1 ± 1.21</b>
Single Arm Throwing	11.5 ± 4.15	12.4 ± 3.11	9.4 ± 2.73	2.5 ± 0.43	<b>16.4 ± 0.99</b>
Single Arm Reorientation	26.2 ± 4.11	1.2 ± 0.58	1.9 ± 0.19	1.1 ± 0.47	<b>36.7 ± 1.81</b>
Two Arm Regrasping	4.1 ± 2.04	0.3 ± 0.15	6.7 ± 1.43	1.7 ± 0.30	<b>24.1 ± 2.55</b>
Two Arm Reorientation	2.2 ± 1.17	0.7 ± 0.45	1.7 ± 1.46	0.5 ± 0.22	<b>35.8 ± 1.64</b>
Two Arm Reorientation (Partition)	1.3 ± 0.94	0.5 ± 0.39	--	--	<b>31.8 ± 6.04</b>
Two Arm Reorientation (Hole)	2.4 ± 0.52	0.5 ± 0.43	--	--	<b>24.7 ± 7.06</b>
Two Arm Reorientation (Barrier)	1.9 ± 0.85	0.4 ± 0.18	--	--	<b>7.7 ± 6.5</b>

Table 1. Mean episode success after  $1e10$  samples ( $5e9$  for the single-arm regrasping task) of different methods. In all tasks, our method outperforms the all baseline methods with relatively small standard deviations.

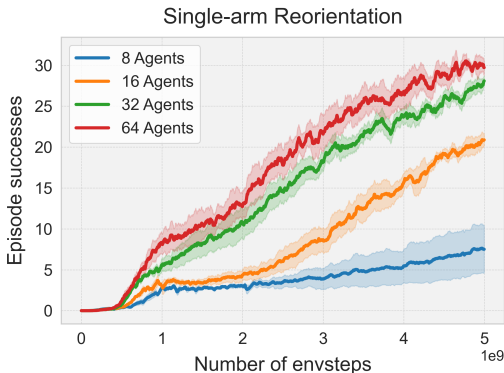


Figure 4. Ablation of agent number ( $K = 8, 16, 32, 64$ ) with the same total number of environments ( $N = 24576$ ). As the number of agents increases, EPO’s performance continuously improves.

EPO is highly robust to hyperparameter variations, such as random seed initialization. This is particularly important given that reinforcement learning algorithms are notoriously sensitive to hyperparameter choices. The robustness of EPO makes it a promising approach for training RL models with minimal hyperparameter tuning, addressing a key challenge in real-world applications.

### 5.3. Ablations

For ablation, we want to see the influence of number of population and GA. We don’t do ablation of hybrid policy update as it would be very similar as pure PBT as shown in section 5.2.

Specifically, we fix the number of environments at  $N = 24,576$  while increasing the population size from  $K = 8$  to  $K = 64$  (Figure 4). The results show that our algorithm benefits from a larger population. While this may seem intuitive, it is not guaranteed by default—as the number of agents increases, behavior diversity also increases, which can lead to more poorly performing agents generating noisy data. By applying Genetic Algorithms (GA) at the latent space level, EPO periodically eliminates underperforming agents, ensuring that only high-quality data is generated

for the master agent to learn from. Additionally, since the actor-critic network is shared across all agents, we introduce controlled diversity while preventing unbounded divergence, which is crucial for scaling.

To further illustrate this, we remove the GA component from EPO, resulting in EPO (w/o GA), which collapses to SAPG. As shown in Figure 5, for harder tasks (e.g., two-arm manipulation), SAPG does not benefit from an increased population at all. For simpler tasks, such as single-arm regrasping, its performance even degrades compared to Figure 3. This empirically demonstrates the necessity of combining GA with hybrid-policy updates to enable effective training at scale with large datasets.

### 5.4. Scaling Law

Here, we examine the scalability of EPO with increased computational resources. Specifically, we evaluate EPO’s scalability with larger training batches generated by parallel simulations by fixing the number of environments per agent and increasing the number of agents from  $K = 16$  to  $K = 128$ , thereby scaling the total number of environments from  $N = 6,144$  to  $N = 49,152$ . It is worth noting that even at  $K = 16$ , the batch size of  $N = 6,144$  is already significantly large. As shown in Figure 6, EPO continues to improve with increasing training data, even when the number of environments reaches  $N = 49,152$ , demonstrating that our algorithm overcomes the limitations of existing on-policy RL methods, which typically saturate beyond a certain batch size. This highlights EPO’s strong potential for data-rich environments, such as simulations and game engines, where large training batches can be easily generated and leveraged effectively.

We also evaluate EPO’s scalability with larger neural networks, where we fix the number of agents and total environments but increase the shared network size. Specifically, we scale the network size by a factor of 4, yet observe no significant increase in per-step training time, while memory usage remains unchanged (because the network used is relatively small). This is not feasible for classic EvoRL methods,



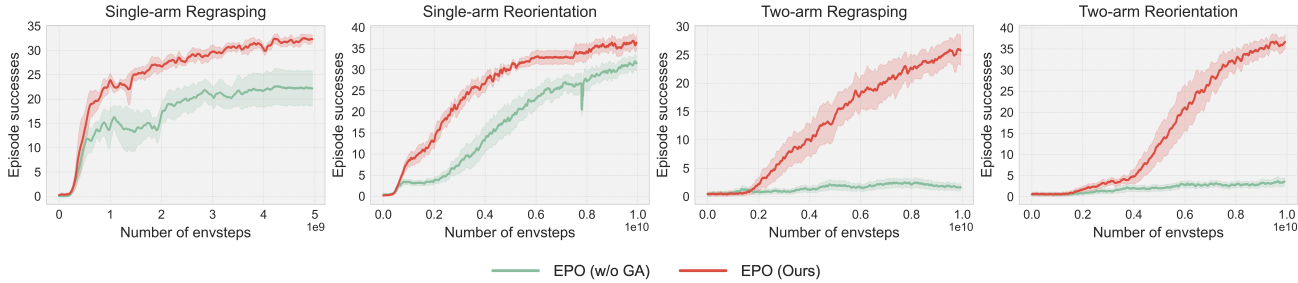


Figure 5. Comparison of the EPO w/ and w/o GA. GA helps EPO maintain the diversity of sampled data while eliminating poorly performing agents, ensuring high-quality data and accelerating training speed.

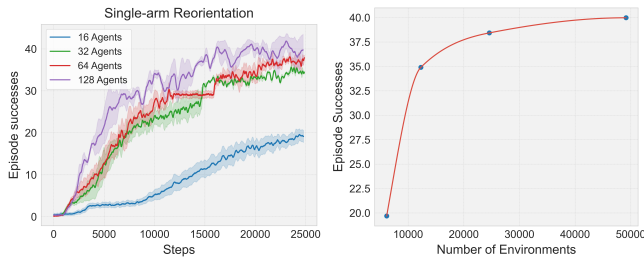


Figure 6. The training curve and scaling curve (converge points) of different numbers of env. With a fixed number of environments per agent (384), as the total number of environments increases, EPO’s performance continues to improve.

where each agent maintains its own separate weights, preventing efficient memory sharing. These results highlight EPO’s strong potential for training large-scale neural networks, including applications in LLMs [12].

## 6. Conclusion

In conclusion, we present Evolutionary Policy Optimization (EPO), a novel policy gradient algorithm that integrates evolutionary algorithms with policy optimization. We demonstrate that EPO significantly outperforms state-of-the-art baselines across several challenging RL benchmarks while also scaling effectively with increased computational resources. Our approach enables large-scale RL training while maintaining high asymptotic performance, and we hope it inspires future research in advancing scalable reinforcement learning.

## Acknowledgements

We would like to thank Ananye Agarwal, Jayesh Singla, Jared Meija, Kevin Gmelin for fruitful discussions.

## Impact Statement

This paper presents Evolutionary Policy Optimization (EPO), a novel reinforcement learning algorithm designed to enhance the scalability and efficiency of policy gradient methods by integrating evolutionary search with on-policy updates. Our work primarily aims to advance the field of Machine Learning (ML) by improving the ability of reinforcement learning (RL) methods to scale with increasing computational resources, making them more effective in data-rich environments such as robotics, simulations, and game engines.

While EPO does not present direct ethical concerns, RL methods, particularly those that improve large-scale training, could have broader societal implications. Scalable RL has potential applications in autonomous systems, decision-making AI, and robotics, which may impact labor markets, automation policies, and safety regulations. It is crucial to ensure that future applications of such methods align with ethical AI principles, particularly regarding fairness, transparency, and robustness in high-stakes decision-making.

Additionally, as RL is increasingly applied to large-scale language models (LLMs) and automated reasoning systems, scalability improvements may accelerate advancements in AI alignment and human-AI interaction. However, misuse of reinforcement learning in generating persuasive AI or manipulative automated systems remains a potential risk. We encourage the community to apply EPO responsibly and consider ethical safeguards when deploying large-scale RL systems in real-world applications.

## References

- [1] Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- [2] Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme,

- S., Hesse, C., et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- [3] Chen, Z., Zhou, Y., He, X., and Jiang, S. A restart-based rank-1 evolution strategy for reinforcement learning. In *IJCAI*, pp. 2130–2136, 2019.
- [4] Choromanski, K., Rowland, M., Sindhvani, V., Turner, R., and Weller, A. Structured evolution with compact architectures for scalable policy optimization. In *International Conference on Machine Learning*, pp. 970–978. PMLR, 2018.
- [5] Choromanski, K. M., Pacchiano, A., Parker-Holder, J., Tang, Y., and Sindhvani, V. From complexity to simplicity: Adaptive es-active subspaces for black-box optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [6] Chrabaszcz, P., Loshchilov, I., and Hutter, F. Back to basics: benchmarking canonical evolution strategies for playing atari. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pp. 1419–1426, 2018.
- [7] Darwin, C. Origin of the species. In *British Politics and the environment in the long nineteenth century*, pp. 47–55. Routledge, 2023.
- [8] Ebert, F., Yang, Y., Schmeckpeper, K., Bucher, B., Georgakis, G., Daniilidis, K., Finn, C., and Levine, S. Bridge data: Boosting generalization of robotic skills with cross-domain datasets. *arXiv preprint arXiv:2109.13396*, 2021.
- [9] Eysenbach, B., Gupta, A., Ibarz, J., and Levine, S. Diversity is all you need: Learning skills without a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- [10] Fujimoto, S., Hoof, H., and Meger, D. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- [11] Gaier, A. and Ha, D. Weight agnostic neural networks. *Advances in neural information processing systems*, 32, 2019.
- [12] Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [13] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018.
- [14] Hoffmeister, F. and Bäck, T. Genetic algorithms and evolution strategies: Similarities and differences. In *International conference on parallel problem solving from nature*, pp. 455–469. Springer, 1990.
- [15] Hwangbo, J., Lee, J., Dosovitskiy, A., Bellicoso, D., Tsounis, V., Koltun, V., and Hutter, M. Learning agile and dynamic motor skills for legged robots. *Science Robotics*, 4(26):eaau5872, 2019.
- [16] Jaderberg, M., Dalibard, V., Osindero, S., Czarnecki, W. M., Donahue, J., Razavi, A., Vinyals, O., Green, T., Dunning, I., Simonyan, K., et al. Population based training of neural networks. *arXiv preprint arXiv:1711.09846*, 2017.
- [17] Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on robot learning*, pp. 651–673. PMLR, 2018.
- [18] Khadka, S. and Tumer, K. Evolution-guided policy gradient in reinforcement learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- [19] Khadka, S., Majumdar, S., Nassar, T., Dwiel, Z., Tumer, E., Miret, S., Liu, Y., and Tumer, K. Collaborative evolutionary reinforcement learning. In *International conference on machine learning*, pp. 3341–3350. PMLR, 2019.
- [20] Kumar, A., Fu, Z., Pathak, D., and Malik, J. Rma: Rapid motor adaptation for legged robots. *Robotics: Science and Systems XVII*, 2021.
- [21] Li, Z., Chen, T., Hong, Z.-W., Ajay, A., and Agrawal, P. Parallel  $q$ -learning: Scaling off-policy reinforcement learning under massively parallel simulation. In *International Conference on Machine Learning*, pp. 19440–19459. PMLR, 2023.
- [22] Lillicrap, T. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [23] Makoviychuk, V., Wawrzyniak, L., Guo, Y., Lu, M., Storey, K., Macklin, M., Hoeller, D., Rudin, N., Allshire, A., Handa, A., et al. Isaac gym: High performance gpu based physics simulation for robot learning. In *Thirty-fifth Conference on Neural Information*

- Processing Systems Datasets and Benchmarks Track (Round 2)*.
- [24] Mania, H., Guy, A., and Recht, B. Simple random search provides a competitive approach to reinforcement learning. *arXiv preprint arXiv:1803.07055*, 2018.
- [25] Meng, W., Zheng, Q., Pan, G., and Yin, Y. Off-policy proximal policy optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 9162–9170, 2023.
- [26] Mnih, V. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [27] Mnih, V. Asynchronous methods for deep reinforcement learning. *arXiv preprint arXiv:1602.01783*, 2016.
- [28] Moriarty, D. E., Schultz, A. C., and Grefenstette, J. J. Evolutionary algorithms for reinforcement learning. *Journal of Artificial Intelligence Research*, 11:241–276, 1999.
- [29] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [30] Petrenko, A., Allshire, A., State, G., Handa, A., and Makovychuk, V. DexPBT: Scaling up dexterous manipulation for hand-arm systems with population based training. *arXiv preprint arXiv:2305.12127*, 2023.
- [31] Pourchot, A. and Sigaud, O. Cem-rl: Combining evolutionary and gradient-based methods for policy search. In *7th International Conference on Learning Representations, ICLR 2019*, 2019.
- [32] Salimans, T., Ho, J., Chen, X., Sidor, S., and Sutskever, I. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- [33] Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust region policy optimization. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1889–1897, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/schulman15.html>.
- [34] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [35] Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., and Riedmiller, M. Deterministic policy gradient algorithms. In *International conference on machine learning*, pp. 387–395. Pmlr, 2014.
- [36] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [37] Singla, J., Agarwal, A., and Pathak, D. Sapg: split and aggregate policy gradients. *arXiv preprint arXiv:2407.20230*, 2024.
- [38] Snoek, J., Larochelle, H., and Adams, R. P. Practical bayesian optimization of machine learning algorithms. *Advances in neural information processing systems*, 25, 2012.
- [39] Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- [40] Tang, Y., Choromanski, K., and Kucukelbir, A. Variance reduction for evolution strategies via structured control variates. In *International Conference on Artificial Intelligence and Statistics*, pp. 646–656. PMLR, 2020.
- [41] Team, G., Anil, R., Borgeaud, S., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., Millican, K., et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [42] Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *nature*, 575(7782):350–354, 2019.
- [43] Zheng, B. and Cheng, R. Rethinking population-assisted off-policy reinforcement learning. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pp. 624–632, 2023.