
Scaling Population-Based Reinforcement Learning with GPU Accelerated Simulation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In recent years, deep reinforcement learning (RL) has shown its effectiveness in
2 solving complex continuous control tasks like locomotion and dexterous manipula-
3 tion. However, this comes at the cost of an enormous amount of experience required
4 for training, exacerbated by the sensitivity of learning efficiency and the policy
5 performance to hyperparameter selection, which often requires numerous trials of
6 time-consuming experiments. This work introduces a Population-Based Reinforce-
7 ment Learning (PBRL) approach that exploits a GPU-accelerated physics simulator
8 to enhance the exploration capabilities of RL by concurrently training multiple
9 policies in parallel. The PBRL framework is applied to three state-of-the-art RL
10 algorithms – PPO, SAC, and DDPG – dynamically adjusting hyperparameters
11 based on the performance of learning agents. The experiments are performed on
12 four challenging tasks in Isaac Gym – *Anymal Terrain*, *Shadow Hand*, *Humanoid*,
13 *Franka Nut Pick* – by analyzing the effect of population size and mutation me-
14 chanisms for hyperparameters. The results demonstrate that PBRL agents outperform
15 non-evolutionary baseline agents across tasks essential for humanoid robots, such
16 as bipedal locomotion, manipulation, and grasping in unstructured environments.
17 The trained agents are finally deployed in the real world for the *Franka Nut Pick*
18 manipulation task. To our knowledge, this is the first sim-to-real attempt for suc-
19 cessfully deploying PBRL agents on real hardware. Code and videos of the learned
20 policies are available on our project website.

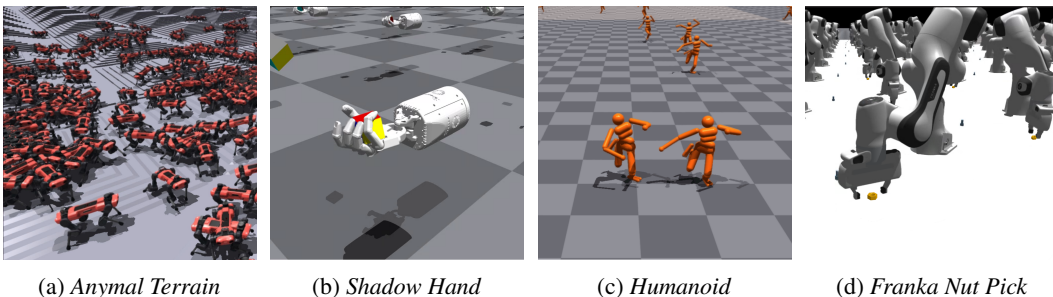


Figure 1: Simulated experiments are performed on four Isaac Gym benchmark tasks: (a) *Anymal Terrain*, to teach a quadruped robot to navigate uneven terrain; (b) *Shadow Hand*, which involves manipulating a cube to a desired orientation with a robot hand; (c) *Humanoid*, for bipedal locomotion; and (d) *Franka Nut Pick*, where the goal is to grasp and lift a nut from a random location on a work surface.

21 1 Introduction

22 Many domains have seen tremendous advancements of reinforcement learning (RL) applications in
23 recent years, ranging from playing challenging games [31, 5] to learning high-dimensional continuous
24 control in robotics [30, 24, 29]. Tasks such as dexterous manipulation [4], legged locomotion [22],
25 and mobile navigation [16] have been learned using deep RL. A primary challenge in training RL
26 policies is the need for large amounts of training data. RL methods rely on effective exploration to
27 discover control policies, which can be particularly challenging when operating in high-dimensional
28 continuous spaces [35]. Moreover, the performance of the learned policy is highly dependent on the
29 tedious tuning of hyperparameters. Hyperparameter tuning can be a very time-consuming process,
30 often requiring many manual trials to determine the best values for the specific task and the learning
31 environment. One way to deal with the problem of data inefficiency is to train in simulation before
32 transferring to reality [32, 2, 23]. However, the time required to train the policy in simulation increases
33 significantly with the task complexity. For example, in [2], learning a block re-orientation task with a
34 robot hand took around 14 days and enormous computing resources. In addition, policies trained
35 in simulation often fail to perform on a real system due to discrepancies between the simulation
36 and the real world. Recent advances in GPU-accelerated simulation, such as Isaac Gym [21, 11],
37 have made it possible to run thousands of parallel environments on a single GPU, which reduces the
38 training times significantly. However, successfully training RL policies still requires carefully tuned
39 hyperparameters to explore efficiently.

40 1.1 Related Works

41 1.1.1 Massively Parallel Simulation

42 The advent of GPU-based simulation has significantly improved simulation throughput by incor-
43 porating massive parallelism on a single GPU [21, 19]. A number of recent works have exploited
44 this parallelism to demonstrate impressive performance on challenging control problems using RL
45 [11, 3, 27]. However, almost all recent works use the same algorithm, i.e. Proximal Policy Optimiza-
46 tion (PPO) [28] to train RL policies; other common approaches include off-policy techniques, e.g.
47 Soft Actor-Critic (SAC) [10] and Deep Deterministic Policy Gradient (DDPG) [20]. While simple
48 and effective, all these algorithms require a range of hyperparameters that need to be tuned for each
49 task to ensure sufficient exploration and stabilize training.

50 1.1.2 Population-Based Reinforcement Learning

51 Population-based approaches offer a promising solution to deal with exploration and hyperparameter
52 tuning by training a set of policies as opposed to a single policy. Multiple agents can be used to
53 collect diverse experiences that improve robustness and stabilize training by dynamically adapting the
54 hyperparameters. Some prior works have shown remarkable results in employing these approaches
55 to train deep RL policies in domains like strategy games and multi-agent interaction [34, 14, 9].
56 However, there is almost no existing research investigating PBRL methods for robotics. This is due
57 to the fact that the computational complexity and training time of these approaches increase linearly
58 with the number of agents on CPU-based simulators like MuJoCo [33], requiring multiple worker
59 machines with separate simulation instances to speed up data collection. Isaac Gym allows simulating
60 thousands of robots in parallel, giving access to a vast amount of experience data, rendering it suitable
61 to efficiently train a population of RL agents.

62 Training various RL agents provides a mechanism for meta-optimization, utilizing the potential of
63 both learning and evolution [1]. One successful example of PBRL methods is population-based
64 training (PBT) [15], which allows training multiple policies concurrently to enhance the exploration
65 capabilities of the agents in generating diverse behaviors. PBT trains a population of agents with
66 different hyperparameters and uses a genetic algorithm to update the population periodically. Recently,
67 DexPBT, a decentralized PBRL approach has been proposed to learn dynamic manipulation between
68 two hand-arm systems using parallel simulations [26]. The authors developed a decentralized
69 implementation to evolve agents in distributed computing environments using on-policy RL, achieving
70 impressive results in dexterous manipulation. However, *sim-to-real* transfer has not been performed,
71 highlighting the complexity of deploying policies on real systems.

72 In contrast, this work targets a broader range of real-world tasks including locomotion and manip-
73 ulation, and transfers the policy onto a real robot without any adaptation phase. In addition, the

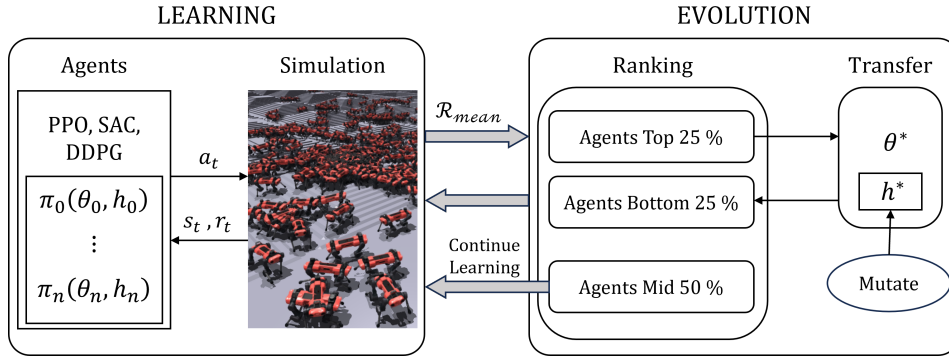


Figure 2: PBRL framework used to learn robotic manipulation tasks through a combination of RL, evolutionary selection, and GPU-based parallel simulations.

74 PBRL framework is successfully applied to both off-policy (SAC, DDPG) and on-policy (PPO) RL
 75 algorithms, analyzing the implications of critical design choices, i.e., the number of agents and the
 76 mutation mechanisms.

77 1.1.3 Sim-to-Real Transfer

78 Despite the calibration efforts to model the physical system accurately, simulation is still a rough
 79 approximation. The differences between the dynamics of simulated and real systems cause a “reality
 80 gap” that makes it unlikely for a simulation-trained policy to successfully transfer to a physical
 81 system. In literature, researchers have put a significant effort into diminishing this gap: to this aim,
 82 most of the approaches leverage domain randomization [24, 4, 3, 27, 8, 7] to expose the policy to
 83 a wide range of observation distributions in simulation, thus improving generalization onto a real
 84 system. Nevertheless, naive domain randomization might not be sufficient to completely attenuate
 85 the dynamics gap: for instance, [13] employs a specific network to mimic the real actuation system.
 86 Another technique in this context is policy-level action integrator (PLAI) [32], a simple yet effective
 87 algorithm aimed at compensating the sim-to-real dynamic discrepancies with an integral action,
 88 which proved to be paramount for a successful transfer.

89 In this paper, we employ sim-to-real strategies to deploy a policy on a real system; to the best of
 90 the authors’ knowledge, this work represents the first instance of deploying PBRL agents on real
 91 hardware.

92 1.2 Contribution

93 This paper investigates a population-based reinforcement learning (PBRL) framework in robotics that
 94 allows the training of a population of agents by exploiting GPU-based massively parallel simulation
 95 to dynamically adjust the hyperparameters during training. We evaluate the PBRL framework on four
 96 complex tasks that require learning essential skills for humanoid robots: *Anymal Terrain*, *Humanoid*,
 97 *Shadow Hand*, and *Franka Nut Pick* (Figure 1), available in Isaac Gym [21]. The results show that
 98 better performance is achieved when training a population of agents compared to a single-agent
 99 baseline on all tasks. The comparison is provided across 3 RL algorithms (PPO, SAC, and DDPG),
 100 varying the number of agents in a population, and across different hyperparameter mutation schemes.
 101 Finally, the PBRL agents are deployed on a real Franka Panda robot for a *Franka Nut Pick* task,
 102 without any policy adaptation phase on the physical system. In summary, the main contributions of
 103 this work are:

- 104 • a population-based RL framework that utilizes GPU-accelerated simulation to train robotic
 105 manipulation tasks by adaptively optimizing the set of hyperparameters during training;
- 106 • simulations demonstrating the effectiveness of the PBRL approach on 4 tasks using 3 RL
 107 algorithms, including both on-policy and off-policy methods, investigating the performance
 108 w.r.t. the number of agents and mutation mechanisms;
- 109 • sim-to-real transfer of PBRL policies onto a real Franka Panda robot;

110 • an open-source codebase to train policies using the PBRL algorithm.

111 2 Methods

112 This section describes the core concepts involved in the PBRL framework. The overall approach,
113 illustrated in Figure 2, can be viewed as a multi-layered training process consisting of an inner
114 optimization loop with RL and an outer loop of online evolutionary selection with population-based
115 training. During training, the parameters of the agent’s policy are updated at a higher rate using RL
116 than the hyperparameters defining the RL procedure.

117 2.1 Reinforcement Learning

118 The RL problem is modeled as a Markov Decision Process (MDP), where an agent interacts with
119 the environment in order to maximize the expected sum of episodic rewards. Specifically, an MDP
120 is defined as $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, \mathcal{T} is the transition
121 dynamics, i.e., $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{P}(\mathcal{S})$, where $\mathbb{P}(\mathcal{S})$ defines the set of a probability distribution over \mathcal{S} ,
122 $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and $\gamma \in [0, 1]$ represents the discount factor. The goal is
123 formulated as learning a policy, either stochastic, $\pi_\theta : \mathcal{S} \rightarrow \mathcal{D}_{\mathcal{A}}$, or deterministic, $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$, where
124 $\mathcal{D}_{\mathcal{A}}$ represents a probability distribution over \mathcal{A} and θ encapsulates the policy parameters, whose
125 cardinality depends on the selected algorithm and network architecture. In this work, the policy is
126 learned using the on-policy method PPO, or either of the off-policy methods SAC or DDPG. All
127 these algorithms use an actor-critic architecture simultaneously learning the policy (actor) and the
128 value function approximators (critics) $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. The implementation of critics in SAC and
129 DDPG relies on double Q-learning and n -step returns.

130 To train the policy with PPO, a learning rate (LR) adaptation procedure is used based on a Kull-
131 back–Leibler (KL) divergence starting from an initial value η_0 [21]. At the end of each update
132 iteration, the LR is increased by a factor of K_η when the KL divergence between the current policy
133 and the old policy is below the specified threshold, or reduced by K_η if the KL divergence exceeds
134 the threshold.

135 In DDPG, the common practice involves adding a small noise to the deterministic actions of the policy
136 to enable exploration. In this work, the noise is added following a mixed exploration strategy [18],
137 where the general idea is akin to adding a different noise level for each environment when training in
138 a massively parallel regime. For the i -th environment out of $N \in \mathbb{Z}^+$ environments, a zero mean and
139 uncorrelated Gaussian noise is given as: $\mathcal{N}(0, \sigma_i)$, where $\sigma_i = \sigma_{min} + \frac{i-1}{N-1}(\sigma_{max} - \sigma_{min})$.

140 2.2 Population-Based Training

141 In standard RL, the agent aims to learn an optimal policy by interacting with an environment and
142 iteratively updating the policy through some kind of optimization method. In contrast, PBRL uses a
143 population of n agents \mathcal{P} , each interacting with the environment independently to collect experience
144 and learn its own policy. Using evolutionary selection, the population is periodically evaluated
145 based on a fitness metric, and best-performing members replace the worst-performing members, i.e.,
146 weights of the best agents are copied over, along with the mutated hyperparameters.

147 In this work, a specific PBRL approach, population-based training (PBT), is employed as an outer
148 optimization loop to enable diverse exploration and dynamically adapt the hyperparameters in high-
149 dimensional continuous control tasks. Each agent $\pi(\theta_i, h_i) \in \mathcal{P}$ is characterized by a vector θ_i and
150 the set of hyperparameters h_i , where θ_i contains the parameters of the policy, and h_i contains the
151 hyperparameters that are optimized during training. To represent the whole population \mathcal{P} , we denote
152 with $\Theta \triangleq \bigcup_{i=1}^n \theta_i$, $\mathbf{h} \triangleq [h_1, h_2, \dots, h_n]$ and $\Pi \triangleq \{\pi(\theta_i, h_i)\}_{i=1}^n$ the sets of all the parameters,
153 hyperparameters and policies respectively.

154 Algorithm 1 provides pseudocode for the PBRL. The training process runs in iterations, where all
155 agents are first independently trained by performing updates to the vector θ_i . After a certain number
156 of policy updates N_{evo} (each agent having been trained for some steps), the agents are evaluated and
157 sorted based on the average return \mathcal{R}_{mean} obtained over all of the previous episodes. The agents in
158 $\mathcal{P}_{bottom\ 25\ \%}$ get replaced by randomly-sampled agents in $\mathcal{P}_{top\ 25\ \%}$ with mutated hyperparameters,
159 while the rest of the agents in $\mathcal{P}_{mid\ 50\ \%}$ and $\mathcal{P}_{top\ 25\ \%}$ continue training.

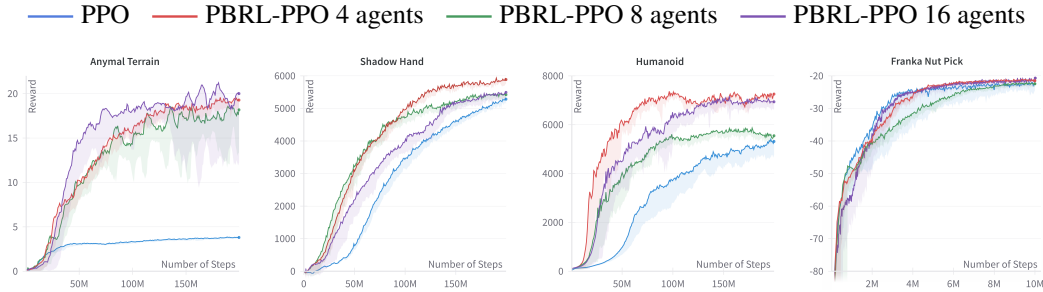


Figure 3: Training results of baseline PPO and PBRL-PPO for $|\mathcal{P}| \in \{4, 8, 16\}$. The shaded area represents the performance between the best and the worst agent in \mathcal{P} , or among 4 different seeds in a non-evolutionary baseline.

160 To generate the mutated hyperparameters, 3 mutation mechanisms are considered (see line 14 of
 161 Algorithm 1): (i) random perturbation is applied to the hyperparameters of the parent agent(s)
 162 through perturbation factors in Table 5; (ii) new hyperparameters are sampled from a prior uniform
 163 distribution with bounds specified in Table 3 and 4; (iii) according to the DexPBT mutation scheme
 164 [26], hyperparameters are multiplied or divided by a random number μ sampled from a uniform
 165 distribution, i.e., $\mu \sim \mathcal{U}(\mu_{min}, \mu_{max})$ with probability $\beta_{mut} \in [0, 1]$. Section 3.2.4 compares all 3
 166 mutation schemes. After beginning the training, evolution is enabled after $N_{start} \in \mathbb{Z}^+$ steps as in
 167 [14] to allow for initial exploration and promote population diversity.

168 3 Experiments

169 3.1 Environments

170 The PBRL framework is evaluated on some of the most challenging benchmark tasks available in
 171 Isaac Gym, including *Anymal Terrain*, *Shadow Hand*, *Humanoid* and *Franka Nut Pick* (Figure 1).
 172 The experiments are conducted on a workstation with a single NVIDIA RTX 4090 GPU and 32 GB
 173 of RAM. Parallelizing the data collection across the GPU, Isaac Gym’s PhysX engine can simulate
 174 thousands of environments using the above hardware.

175 3.2 Results

176 The experiments focus on optimizing the hyperparameters of the RL agents in a population and
 177 comparing the results against non-evolutionary baseline agents. For each case of baseline agents,
 178 4 experiments are run with different seeds. Tables 3 and 4 list the hyperparameters for on-policy
 179 and off-policy algorithms, including the sampling ranges of those optimized through the PBRL
 180 Algorithm 1. The initial values for each agent are uniformly sampled from a prior distribution with a
 181 given range.

182 3.2.1 PBRL-PPO

183 For the PPO agents, the tuned hyperparameters are the KL divergence threshold for an adaptive
 184 LR, the entropy loss coefficient, and the variance of action selection. These parameters are crucial
 185 in ensuring sufficient exploration of the environment. Figure 3 shows the learning curves for the
 186 single-agent PPO baseline and PBRL-PPO for $|\mathcal{P}| \in \{4, 8, 16\}$. The results demonstrate that PBRL-
 187 PPO outperforms PPO on 3 out of 4 tasks, yielding a higher return, with significant improvement
 188 seen in *Anymal Terrain*, which involves traversing increasingly challenging terrain. For *Franka Nut*
 189 *Pick*, PBRL agents achieve comparable performance to the baseline PPO agents. This is because, in
 190 this relatively straightforward task, randomization alone suffices for a thorough exploration of the
 191 state/action space.

192 3.2.2 PBRL-SAC

193 In PBRL-SAC, the optimized hyperparameters include the LR of the actor-critic networks and the
 194 target entropy factor. Entropy is key in SAC agents as the policy is trained to maximize the trade-

195 off between the expected return and exploration. Experiments are run with a population size of
 196 $|\mathcal{P}| \in \{4, 8\}$. Due to higher memory needs for replay buffers in off-policy methods, the maximum
 197 population size is limited to 8. The training performance of SAC and PBRL-SAC is shown in
 198 Figure 4. PBRL-SAC improves the training performance compared to non-evolutionary SAC on 3
 199 out of 4 tasks, yielding a remarkable improvement on both *Shadow Hand* and *Franka Nut Pick*, while
 200 comparable results are achieved on *Humanoid*, probably due to the limited task complexity.

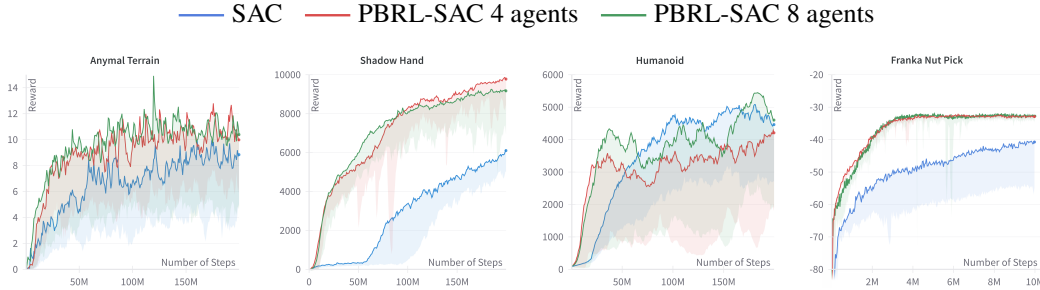


Figure 4: Training results of baseline SAC and PBRL-SAC for $|\mathcal{P}| \in \{4, 8\}$. The shaded area displays the performance between the best and the worst agent in \mathcal{P} , or among 4 different seeds in a non-evolutionary baseline.

201 3.2.3 PBRL-DDPG

202 In DDPG, exploration noise is added to the output of the deterministic actor. As mentioned in
 203 Section 2.1, different noise levels are added for different environments uniformly within the range
 204 $[\sigma_{min}, \sigma_{max}]$. Both these parameters are crucial in controlling the amount of exploration in DDPG
 205 agents. In PBRL-DDPG, the hyperparameters optimized during training include the minimum and
 206 the maximum bounds for noise levels, i.e., $\sigma_{min}, \sigma_{max}$, and the LRs of the actor and the critic. As
 207 in PBRL-SAC, the maximum population size in PBRL-DDPG is set to 8 due to the presence of
 208 independent replay buffers and GPU memory limitations. Figure 5 shows that PBRL-DDPG achieves
 209 significantly better training performance than DDPG on all 4 benchmark tasks.

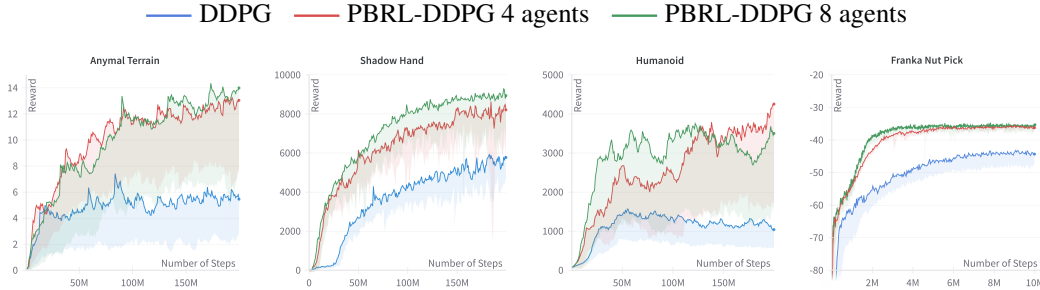


Figure 5: Training results of baseline DDPG and PBRL-DDPG for $|\mathcal{P}| \in \{4, 8\}$. The shaded area displays the performance between the best and the worst agent in \mathcal{P} , or among 4 different seeds in a non-evolutionary baseline.

210 3.2.4 Mutation Comparison

211 Figure 6 shows the results using 3 different mutation schemes for PBRL-PPO and PBRL-DDPG. As
 212 mentioned in Section 2.2, the hyperparameters for under-performing agents are generated either by
 213 sampling from an original prior distribution, by perturbing the parent’s values through perturbation
 214 factors given in Table 5, or through the DexPBT mutation scheme presented in [26]. In the latter,
 215 the hyperparameters have a $\beta_{mut} := 0.5$ probability of getting multiplied or divided by a random
 216 number sampled from the uniform distribution, $\mu \sim \mathcal{U}(1.1, 1.5)$. The results show that the perturbed
 217 agents either exceed or are on par with the performance of other mutation schemes in 6 out of 8
 218 experiments. The DexPBT mutation scheme performs better with PBRL-DDPG on *Humanoid* and
 219 *Franka Nut Pick* tasks, which are less challenging compared to others. The combination of two

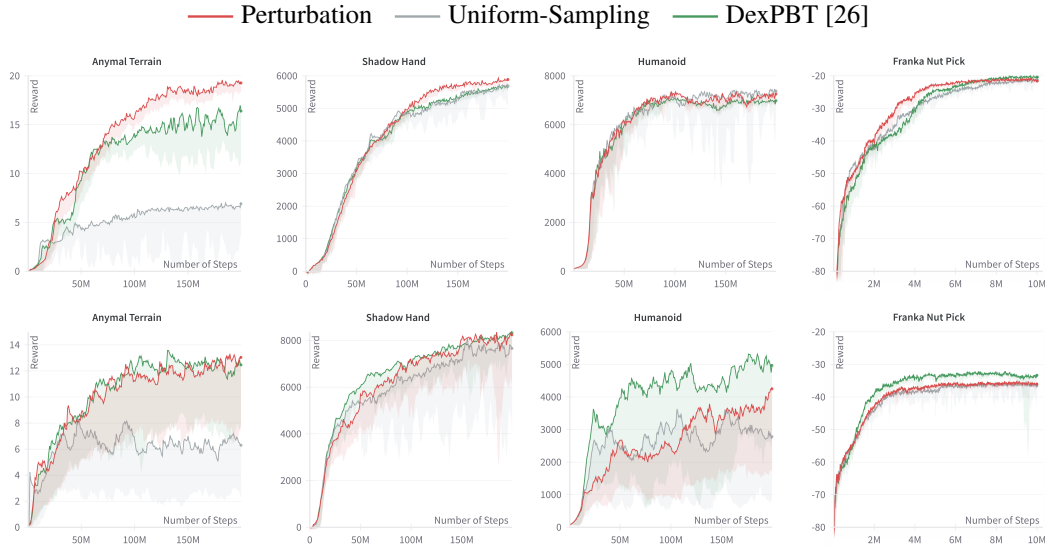


Figure 6: Comparison of different mutation schemes for PBRL-PPO (top) and PBRL-DDPG (bottom) with $|\mathcal{P}| = 4$.

220 mutation schemes might discover better exploration strategies for a wider range of tasks. Analyzing
 221 the potential synergies between the two remains a prospect for future investigation.

222 3.3 Sim-to-Real Transfer

223 In the real experiments, we replicate the *Franka Nut Pick* task [24] by deploying a PBRL-PPO policy,
 224 without any real-world adaptation, executing the actions with PLAI [32]. The robot detects the nuts
 225 utilizing Mask-RCNN [12], fine-tuned on real-world images captured with a wrist-mounted Intel
 226 RealSense D435 camera, using the `IndustRealLib` codebase [32]. Compared to the original task
 227 [24], we applied the following changes to make the simulated environment resemble real setup:
 228 (i) employing a Task-Space Impedance (TSI) controller [6] instead of an Operational-Space motion
 229 Controller (OSC) [17] to comply with the actual low-level controller¹; (ii) randomizing the nut’s
 230 initial position to reflect the actual robot workspace; (iii) changing the observation space to include
 231 the 7-DOF joint configuration, the measured end-effector pose, and the estimated nut pose. The
 232 parameters used in the simulated environment and the real controller are reported in Table 2.

233 During experiments, the following policies were deployed, performing 30 real-world trials of *Franka*
 234 *Nut Pick* task for each policy: (i) 2 agents from a population of 8, trained with PBRL-PPO, specifically
 235 the “best” and the “worst” agent; (ii) the “best” agent trained with baseline PPO. With “best” and
 236 “worst” we indicate the agents achieving the highest and lowest *success rate* in simulation, where
 237 success is defined as reaching, grasping, and lifting the nut, without losing contact during the lifting
 238 phase. PBRL-PPO with $|\mathcal{P}| = 8$ achieved the highest success rate. Remarkably, we found out that
 239 even the success rate of the worst agent in \mathcal{P} was higher than that of the best PPO agent.

240 Deploying both PPO and PBRL-PPO agents onto a real robot leads to task completion (shown in
 241 Figure 7), yet with different success rates, as summarized in Table 1. Particularly, both PBRL-PPO

¹The control laws are specified in [24] and in reference works [6, 17]

Table 1: Success rate deploying the best and the worst of 8 agents trained with PBRL-PPO and the best PPO baseline agent on the *Franka Nut Pick* task with the real robot

Algorithm	Agent	Successful trials	Success rate
PBRL-PPO	Best	27/30	90%
PBRL-PPO	Worst	21/30	70%
PPO	Best	19/30	63.33%

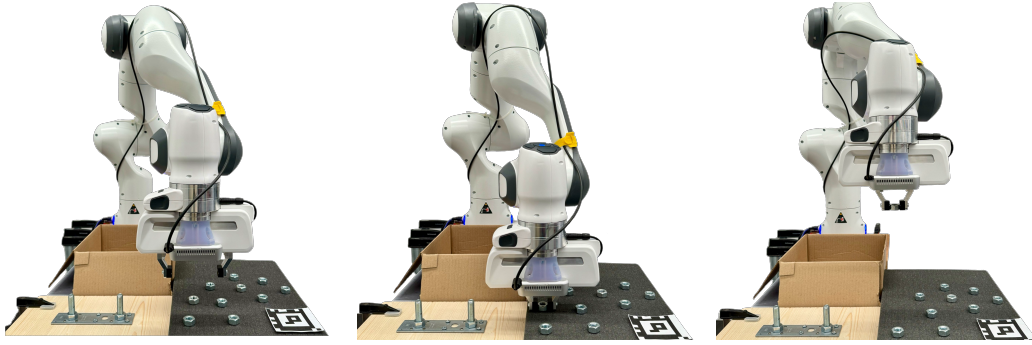


Figure 7: Snapshots of the *Franka Nut Pick* experiment on the real robot: full video on our project website.

242 agents yield higher success rates than PPO, with the “best” agent performing better than the “worst”
243 one, indeed confirming the ranking attained in simulation. Unlike the baseline PPO agent, which
244 continued to produce small movements after reaching the target, PBRL-PPO agents remained more
245 stable, leading to a higher success rate. This demonstrates that PBRL agents, while achieving similar
246 rewards to a single agent, learn behaviors that exhibit greater robustness to environment variability
247 due to the diversity in agent populations. Informally, the best PBRL-PPO agent also exhibited
248 recovery behavior during task execution after perturbation by the human.

249 3.4 Discussion

250 While the PBRL agents perform better than the non-evolutionary agents in almost all the experiments,
251 the impact of population size across RL algorithms and tasks shows no consistent pattern. One
252 may hypothesize that larger and more diverse populations might lead to a better final performance.
253 However, the results in this work indicate that using a larger population size does not necessarily yield
254 substantial benefits for every task. This is in contrast to the common belief that population-based
255 methods rely on larger population sizes to effectively explore the hyperparameter space [15, 25]. The
256 optimal population size, instead, depends on various factors, including task complexity, RL algorithm,
257 and interaction dynamics among agents. While larger populations offer increased exploration
258 potential, they also suffer from diminished exploitation capabilities due to increased competition,
259 leading to lower performance in less challenging tasks where smaller populations suffice. Larger
260 population sizes seem to perform better when the task complexity gradually increases requiring
261 extensive exploration as in *Anymal Terrain*, which implements curriculum learning.

262 Additionally, the performance of PBRL may be lower than non-evolutionary agents on relatively
263 simpler tasks where optimal hyperparameters are known *a priori*. This can be noticed on a *Humanoid*
264 task trained with SAC in Figure 4: indeed, baseline policies achieve a higher reward than PBRL-SAC
265 with 4 agents; nevertheless, 8 agents are capable of outperforming the baseline. Thus, the benefits
266 provided by PBRL will become more apparent for new tasks where ideal hyperparameter ranges are
267 not known in advance. In this sense, PBRL can be thought of as an exploratory approach to search
268 for unknown optimal configurations of newly designed tasks.

269 4 Conclusion

270 In this paper, a PBRL framework has been employed to train a population of RL agents by making use
271 of high-throughput GPU-accelerated simulation. The first simulation results of PBRL using on-policy
272 and off-policy methods are provided on a series of locomotion and manipulation benchmark tasks
273 proposed in [21] by investigating the effect of population size and different mutation schemes. The
274 results showed the effectiveness of PBRL in improving final performance through online adaptation of
275 hyperparameters. PBRL agents have been deployed on real hardware for the first time, demonstrating
276 smooth and successful transfer, without any policy adaptation or fine-tuning. Finally, we released
277 the codebase to train PBRL agents and hope that it will empower researchers to further explore and
278 extend the capabilities of PBRL algorithms.

References

- 279
280 [1] D. Ackley and M. Littman. Interactions between learning and evolution. *Artif. Life II*, 10:487–
281 509, 1991.
- 282 [2] I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino,
283 M. Plappert, G. Powell, R. Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint:*
284 *1910.07113*, 2019.
- 285 [3] A. Allshire, M. Mittal, V. Lodaya, V. Makoviychuk, D. Makoviichuk, F. Widmaier, M. Wüthrich,
286 S. Bauer, A. Handa, and A. Garg. Transferring dexterous manipulation from gpu simulation to
287 a remote real-world trifinger. In *Proc. IEEE Int. Conf. Intell. Robots Syst.*, pages 11802–11809,
288 Oct. 2022.
- 289 [4] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron,
290 M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *Int. J. Robot.*
291 *Res.*, 39(1):3–20, Jan. 2020.
- 292 [5] C. Berner, G. Brockman, B. Chan, V. Cheung, P. Debiak, C. Dennison, D. Farhi, Q. Fischer,
293 S. Hashme, C. Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint:*
294 *1912.06680*, 2019.
- 295 [6] F. Caccavale, C. Natale, B. Siciliano, and L. Villani. Six-DOF impedance control based on
296 angle/axis representations. *IEEE Trans. Robot. Automat.*, 15(2):289–300, Apr. 1999.
- 297 [7] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox. Closing
298 the Sim-to-Real Loop: Adapting Simulation Randomization with Real World Experience. In
299 *Proc. IEEE Int. Conf. Robot. Automat.*, pages 8973–8979, May 2019.
- 300 [8] C. Chi, B. Burchfiel, E. Cousineau, S. Feng, and S. Song. Iterative Residual Policy for Goal-
301 Conditioned Dynamic Manipulation of Deformable Objects. In *Proc. Robot. Sci. Syst.*, June
302 2022.
- 303 [9] A. Flajolet, C. B. Monroc, K. Beguir, and T. Pierrot. Fast population-based reinforcement
304 learning on a single machine. In *Proc. Int. Conf. Mach. Learn.*, volume 162, pages 6533–6547,
305 July 2022.
- 306 [10] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy
307 deep reinforcement learning with a stochastic actor. *arXiv preprint: 1801.01290*, 2018.
- 308 [11] A. Handa, A. Allshire, V. Makoviychuk, A. Petrenko, R. Singh, J. Liu, D. Makoviichuk,
309 K. Van Wyk, A. Zhurkevich, B. Sundaralingam, and Y. Narang. Dextreme: Transfer of agile
310 in-hand manipulation from simulation to reality. In *Proc. IEEE Int. Conf. Robot. Automat.*,
311 pages 5977–5984, June 2023.
- 312 [12] K. He, G. Gkioxari, P. Dollar, and R. Girshick. Mask R-CNN. In *Proc. IEEE Int. Conf. Comput.*
313 *Vis.*, pages 2980–2988, Oct. 2017.
- 314 [13] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter. Learning
315 agile and dynamic motor skills for legged robots. *Sci. Robot.*, 4(26), Jan. 2019.
- 316 [14] M. Jaderberg, W. M. Czarnecki, I. Dunning, L. Marris, G. Lever, A. G. Castañeda, C. Beattie,
317 N. C. Rabinowitz, A. S. Morcos, A. Ruderman, N. Sonnerat, T. Green, L. Deason, J. Z. Leibo,
318 D. Silver, D. Hassabis, K. Kavukcuoglu, and T. Graepel. Human-level performance in 3D
319 multiplayer games with population-based reinforcement learning. *Science*, 364(6443):859–865,
320 May 2019.
- 321 [15] M. Jaderberg, V. Dalibard, S. Osindero, W. M. Czarnecki, J. Donahue, A. Razavi, O. Vinyals,
322 T. Green, I. Dunning, K. Simonyan, et al. Population based training of neural networks. *arXiv*
323 *preprint*, 2017.
- 324 [16] G. Kahn, A. Villafior, B. Ding, P. Abbeel, and S. Levine. Self-supervised deep reinforcement
325 learning with generalized computation graphs for robot navigation. In *Proc. IEEE Int. Conf.*
326 *Robot. Automat.*, pages 5129–5136, May 2018.

- 327 [17] O. Khatib. A unified approach for motion and force control of robot manipulators: The
328 operational space formulation. *IEEE J. Robot. Automat.*, 3(1):43–53, Feb. 1987.
- 329 [18] Z. Li, T. Chen, Z.-W. Hong, A. Ajay, and P. Agrawal. Parallel Q -learning: Scaling off-policy
330 reinforcement learning under massively parallel simulation. In *Proc. Int. Conf. Mach. Learn.*,
331 volume 202, pages 19440–19459, July 2023.
- 332 [19] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox. GPU-accelerated
333 robotic simulation for distributed reinforcement learning. In *Proc. Conf. Robot Learn.*, vol-
334 ume 87, pages 270–282, Oct. 2018.
- 335 [20] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.
336 Continuous control with deep reinforcement learning. *arXiv preprint: 1509.02971*, 2015.
- 337 [21] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin,
338 A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for
339 robot learning. *arXiv preprint: 2108.10470*, 2021.
- 340 [22] G. Margolis, G. Yang, K. Paigwar, T. Chen, and P. Agrawal. Rapid Locomotion via Reinforce-
341 ment Learning. In *Proc. Robot. Sci. Syst.*, June 2022.
- 342 [23] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter. Learning robust
343 perceptive locomotion for quadrupedal robots in the wild. *Sci. Robot.*, 7(62), Jan. 2022. Art. no.
344 eabk2822.
- 345 [24] Y. Narang, K. Storey, I. Akinola, M. Macklin, P. Reist, L. Wawrzyniak, Y. Guo, A. Moravanszky,
346 G. State, M. Lu, A. Handa, and D. Fox. Factory: Fast Contact for Robotic Assembly. In *Proc.*
347 *Robot. Sci. Syst.*, June 2022.
- 348 [25] J. Parker-Holder, V. Nguyen, and S. J. Roberts. Provably efficient online hyperparameter
349 optimization with population-based bandits. In *Proc. Adv. Neural Inform. Process. Syst.*,
350 volume 33, pages 17200–17211, Dec. 2020.
- 351 [26] A. Petrenko, A. Allshire, G. State, A. Handa, and V. Makoviychuk. DexPBT: Scaling up
352 Dexterous Manipulation for Hand-Arm Systems with Population Based Training. In *Proc.*
353 *Robot. Sci. Syst.*, July 2023.
- 354 [27] N. Rudin, D. Hoeller, P. Reist, and M. Hutter. Learning to walk in minutes using massively
355 parallel deep reinforcement learning. In *Proc. Conf. Robot Learn.*, volume 164, pages 91–100,
356 Nov. 2022.
- 357 [28] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
358 algorithms. *arXiv preprint: 1707.06347*, 2017.
- 359 [29] A. A. Shahid, D. Piga, F. Braghin, and L. Roveda. Continuous control actions learning and
360 adaptation for robotic manipulation through reinforcement learning. *Auton. Robots*, 46(3):483–
361 498, Feb. 2022.
- 362 [30] A. A. Shahid, J. S. V. Sesin, D. Pecioski, F. Braghin, D. Piga, and L. Roveda. Decentralized
363 multi-agent control of a manipulator in continuous task learning. *Appl. Sci.*, 11(21):Art. no.
364 10227, Nov. 2021.
- 365 [31] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre,
366 D. Kumaran, T. Graepel, et al. A general reinforcement learning algorithm that masters chess,
367 shogi, and go through self-play. *Science*, 362(6419):1140–1144, Dec. 2018.
- 368 [32] B. Tang, M. A. Lin, I. A. Akinola, A. Handa, G. S. Sukhatme, F. Ramos, D. Fox, and Y. S.
369 Narang. IndustReal: Transferring Contact-Rich Assembly Tasks from Simulation to Reality. In
370 *Proc. Robot. Sci. Syst.*, July 2023.
- 371 [33] E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In *Proc.*
372 *IEEE Int. Conf. Intell. Robots Syst.*, pages 5026–5033, Oct. 2012.

- 373 [34] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi,
374 R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster level in StarCraft II using multi-agent
375 reinforcement learning. *Nature*, 575(7782):350–354, Nov. 2019.
- 376 [35] J. Xu, V. Makoviychuk, Y. Narang, F. Ramos, W. Matusik, A. Garg, and M. Macklin. Accelerated
377 policy learning with parallel differentiable simulation. In *Proc. Int. Conf. Learn. Represent.*,
378 Apr. 2022. Art. no. 186704.

379 **A Algorithm**

Algorithm 1 PBRL algorithm

Require: Initial population \mathcal{P} of agents (Θ random, h sampled from a uniform distribution)

- 1: $N_{iter} = 0$
- 2: **while** not end of training **do**
- 3: $\theta \leftarrow \text{Train}(\Pi(\Theta, h))$ ▷ Train all agents in \mathcal{P}
- 4: $N_{iter} = N_{iter} + 1$
- 5: **if** $N_{iter} > N_{start}$ and $N_{iter} \% N_{evo} = 0$ **then**
- 6: **for** each agent $\pi(\theta, h) \in \mathcal{P}$ **do**
- 7: $\mathcal{R}_{mean} \leftarrow \text{Eval}(\pi(\theta, h))$
- 8: Sort $\pi(\theta, h)$ based on \mathcal{R}_{mean}
- 9: **end for**
- 10: Partition \mathcal{P} into $\mathcal{P}_{top\ 25\%}$, $\mathcal{P}_{mid\ 50\%}$, $\mathcal{P}_{bottom\ 25\%}$
- 11: Sample $\pi^*(\theta^*, h^*)$ from $\mathcal{P}_{top\ 25\%}$ at random
- 12: **for** each agent $\pi(\theta, h) \in \mathcal{P}_{bottom\ 25\%}$ **do**
- 13: $\pi(\theta) \leftarrow \pi^*(\theta^*)$
- 14: $h \leftarrow \text{Mutate}(h^*)$
- 15: **end for**
- 16: **end if**
- 17: **end while**

380 **B Domain Randomization for Franka Nut Pick Task**

381 In this section, we include the settings used for domain randomization in experiments with the Franka
 382 robotic arm for simulated environment and real setup. The robot initial pose is randomized according
 to a Gaussian distribution \mathcal{N} , while the nut initial position is uniformly chosen in the specified range.

Table 2: Simulated environment and real control configuration parameters used in *Franka Nut Pick* during training and deployment respectively.

Parameter	Value
Franka initial position	$\mathcal{N}([0.0, -0.2, 0.2], [0.2, 0.2, 0.1])$
Franka initial rotation	$\mathcal{N}([\pi, 0, \pi], [0.3, 0.3, 1])$
Nut initial position	$[0.42, 0.27, 0.02] \pm [0.18, 0.13, 0.01]$
TSI proportional gains	$[1000, 1000, 1000, 50, 50, 50]$
TSI derivative gains	$[63.25, 63.25, 63.25, 1.414, 1.414, 1.414]$
Action scale	0.0001

383

384 **C Hyperparameters**

Table 3: Hyperparameters setup for PPO and PBRL-PPO across all the tasks.

Hyperparameter	PPO			PBRL-PPO		
	<i>Anymal Terrain</i>	<i>Shadow Hand & Humanoid</i>	<i>Franka Nut Pick</i>	<i>Anymal Terrain</i>	<i>Shadow Hand & Humanoid</i>	<i>Franka Nut Pick</i>
Environments per agent	4096	16384	128	1024	4096	128
MLP hidden units	[512, 256, 128]	[512, 256, 128]	[256, 128, 64]	[512, 256, 128]	[512, 256, 128]	[256, 128, 64]
Horizon	32	16	120	32	16	120
Batch size	8192	32768	512	8192	8192	512
Actor variance	0.5	1	1	0.3 - 1	0.3 - 1	0.3 - 1
KL threshold	0.016	0.016	0.016	0.08 - 0.016	0.08 - 0.016	0.08 - 0.016
Entropy loss coefficient	0.001	0.001	0	0.0001 - 0.001	0.0001 - 0.001	0.0001 - 0.001
Epochs	8	4	8	8	4	8
Discount factor γ	0.99	0.99	0.99	0.99	0.99	0.99
GAE lambda	0.95	0.95	0.95	0.95	0.95	0.95
PPO clip ϵ	0.2	0.2	0.2	0.2	0.2	0.2
Initial LR η_0	5×10^{-4}	5×10^{-4}	5×10^{-4}	5×10^{-4}	5×10^{-4}	5×10^{-4}
LR adaptation gain K_η	1.5	1.5	1.5	1.5	1.5	1.5

Table 4: Hyperparameters setup for off-policy algorithms on all four tasks. *For *Franka Nut Pick* these parameters are, respectively: 128, [256, 128, 64], 512.

Hyperparameter	SAC & DDPG	PBRL-SAC & PBRL-DDPG
Environments per agent*	2048	2048
MLP hidden units*	[512, 256, 128]	[512, 256, 128]
Batch size*	4096	4096
Horizon	1	1
Target update rate τ	5×10^{-2}	5×10^{-2}
Actor learning rate	0.0001	0.0001 – 0.001
Critic learning rate	0.0001	0.0001 – 0.001
DDPG exploration σ_{min}	0.01	0.01 – 0.1
DDPG exploration σ_{max}	1	0.5 – 1
SAC target entropy	-20	-20 – -10
Replay buffer size	1×10^6	1×10^6
Epochs	4	4
n -step returns	3	3

Table 5: Parameter setup for PBRL

Parameter	Value	
	<i>Franka Nut Pick</i>	Others
Evolution start N_{start}	2×10^5 steps	1×10^7 steps
Evolution frequency N_{evo}	1×10^5 steps	2×10^6 steps
Perturbation factor (min.)	0.8	0.8
Perturbation factor (max.)	1.2	1.2