
MULTIPLE-FREQUENCIES POPULATION-BASED TRAINING

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement Learning’s high sensitivity to hyperparameters is a source of instability and inefficiency, creating significant challenges for practitioners. Hyperparameter Optimization (HPO) algorithms have been developed to address this issue, among them Population-Based Training (PBT) stands out for its ability to generate hyperparameters schedules instead of fixed configurations. PBT trains a population of agents, each with its own hyperparameters, frequently ranking them and replacing the worst performers with mutations of the best agents. These intermediate selection steps can cause PBT to focus on short-term improvements, leading it to get stuck in local optima and eventually fall behind vanilla Random Search over longer timescales. This paper studies how this greediness issue is connected to the choice of *evolution frequency*, the rate at which the selection is done. We propose Multiple-Frequencies Population-Based Training (MF-PBT), a novel HPO algorithm that addresses greediness by employing sub-populations, each evolving at distinct frequencies. MF-PBT introduces a migration process to transfer information between sub-populations, with an asymmetric design to balance short and long-term optimization. Extensive experiments on the Brax suite demonstrate that MF-PBT improves sample efficiency and long-term performance, even without tuning hyperparameters. Code will be released.

1 INTRODUCTION

The performance of neural networks depends on selecting a well-suited configuration of hyperparameters, a task that is time-consuming and often reduced to trial-and-error when done manually. This concern has driven the development of Hyperparameter Optimization (HPO, Bergstra et al. (2011); Feurer & Hutter (2019)), field focused on modeling and automating the hyperparameter selection process. The need for HPO algorithms is particularly strong in Reinforcement Learning (RL, Sutton & Barto (2018)), as RL algorithms are often highly sensitive to hyperparameter choices (Eimer et al., 2023; Zhang et al., 2021).

Given these challenges, Population-Based Training (PBT, Jaderberg et al. (2017)) has become a popular HPO method among RL practitioners (Badia et al., 2020; Liu et al., 2022). PBT trains a population of agents in parallel, using an evolutionary process to propagate successful hyperparameter configurations while exploring new ones. This frequent evolution enables PBT to generate dynamic hyperparameter schedules, unlike earlier methods like random search (Bergstra & Bengio, 2012) and classic sequential optimization (Li et al., 2018; Falkner et al., 2018), which typically produced fixed configurations. This dynamic adaptation of hyperparameters is particularly desirable in RL, where the learning problem is non-stationary (Parker-Holder et al., 2022).

However, to achieve this dynamic adaptation, PBT selects hyperparameter configurations based on intermediate performance. As a result, it often favors configurations that show early improvements but fail to deliver better long-term results. Dalibard & Jaderberg (2021) identified this greediness and proposed Faster Improvement Rate PBT (FIRE PBT), which uses learning curves to predict the long-term potential of hyperparameters based on their improvement rates. In this paper, we address PBT’s inherent greediness by introducing a novel focus on evolution frequencies.

Evolution frequency, which controls the number of training steps between evolutionary updates, has not been explicitly addressed in prior research on PBT. Yet, our study shows that it lies at the core of a key trade-off in PBT’s behavior. Evolving too frequently can lead to greedy collapse in two ways: (1)

054 aggressive hyperparameter tuning traps PBT in local optima, and (2) population diversity decreases as
055 similar agents are reproduced repeatedly. Conversely, reducing the evolution frequency limits PBT’s
056 adaptability, resulting in less fine-grained schedules and, ultimately, deteriorating sample efficiency.

057 To address this trade-off, we propose Multiple-Frequencies Population-Based Training (MF-PBT),
058 a novel HPO algorithm that employs multiple sub-populations, each evolving at distinct frequen-
059 cies. By incorporating an asymmetric migration process, MF-PBT allows these sub-populations to
060 share information while preventing greediness. This design aims to balance short and long-term
061 optimization, leading to mixed-frequency schedules that enhance anytime performance.

062 We validate MF-PBT through a series of reinforcement learning experiments using the Brax frame-
063 work (Freeman et al., 2021). Our results demonstrate that MF-PBT effectively mitigates the two
064 forms of greedy collapse, achieving significantly higher long-term rewards and improved anytime
065 performance compared to PBT baselines. Additionally, we conduct an empirical study on the potential
066 of population-based methods for *variance-exploitation*, showing that even without hyperparameter
067 tuning, populations can greatly improve performance by exploiting the inherent stochasticity of RL
068 training. To ensure reproducibility, we make our code publicly available.

069 To summarize, our contributions are as follows:
070

- 071 1. We investigate the impact of evolution frequency on PBT and its connection to greediness.
- 072 2. We introduce MF-PBT, a novel HPO algorithm that uses multiple evolution frequencies and
073 an asymmetric migration process across sub-populations to overcome PBT’s greediness
- 074 3. We evaluate MF-PBT using the Brax suite, isolating the impact of PBT’s greediness and
075 demonstrating how MF-PBT mitigates this issue to achieve better final performance and
076 sample efficiency across various environments.
- 077 4. We empirically show how population-based methods can leverage stochasticity in rein-
078 forcement learning training to significantly improve performance, even without explicit
079 hyperparameter tuning.

082 2 PRELIMINARIES

084 Hyperparameter Optimization (HPO, Bergstra et al. (2011); Feurer & Hutter (2019)) encompasses
085 various approaches aimed at efficiently tuning hyperparameters to enhance performance and robust-
086 ness of learning algorithms. Random Search (Bergstra & Bengio, 2012) is the baseline approach
087 to HPO; the field then progressed towards more sophisticated techniques, notably meta-gradient
088 methods (Finn et al., 2017; Xu et al., 2018), sequential optimization (Li et al., 2018; Falkner et al.,
089 2018; Awad et al., 2021), and population-based approaches.

090 In this section, we focus on Population-Based Training (PBT, Jaderberg et al. (2017)), beginning with
091 its mechanisms and relevance to RL applications. Next, we briefly review notable extensions to PBT
092 and provide insights into the greediness issue that we aim to tackle, highlighting its connection to
093 evolution frequency.

095 2.1 POPULATION-BASED TRAINING

097 Population-Based Training (PBT) is an HPO technique that combines evolutionary strategies with
098 gradient-based optimization. In PBT, a population of N agents, $\mathcal{P} = \{a_i\}_{i=1}^N$, is trained iteratively
099 in parallel, with each agent maintaining its own set of hyperparameters, h_i , and neural network
100 parameters, θ_i .

101 After every t_{ready} training steps, an *evolution step* occurs where all agents are evaluated and assigned
102 a fitness score. The parameter t_{ready} controls the *evolution frequency*, with smaller values resulting
103 in more frequent evolution. The agents are then ranked and divided into three brackets: *winner*s,
104 *survivor*s, and *loser*s. The evolution step consists of two phases: an exploitation phase, where the
105 losers are replaced with clones of the winners, followed by an exploration phase, where the cloned
106 agents’ hyperparameters are slightly perturbed to encourage exploration around the best solutions.
107 In our experiments, we use the *truncation* method introduced in PBT: the top 25% are winners, the
bottom 25% are losers, and the remaining agents are survivors.

108 While PBT can be applied to any deep learning task, it is particularly effective in RL, due to the
109 non-stationarity of the training process (Parker-Holder et al., 2022; Zhang et al., 2021). Unlike
110 supervised learning, where the data distribution remains fixed, RL experiences significant shifts in
111 the data distribution as training progresses, and the hyperparameters should take it into account.
112 PBT’s frequent evolution steps allow hyperparameters to adapt to the current learning state, naturally
113 generating schedules that accommodate this non-stationarity.

114 Another strength of PBT is its ability to harness RL’s intrinsic variance. The stochastic nature of
115 both the environment and learning algorithms leads to significant performance fluctuations across
116 different random seeds (Henderson et al., 2018; Agarwal et al., 2021). By maintaining a population
117 and periodically reproducing the top performers, PBT propagates favorable outcomes, ensuring
118 that unfortunate agents are replaced by luckier ones. This ability of PBT to propagate exploration
119 luck is noted in Jaderberg et al. (2017), but our experiments in section 4.3 further demonstrate that
120 population-based approaches can significantly improve performance, even without hyperparameter
121 tuning.

122 Numerous extensions to PBT have been proposed, focusing on improving exploration and efficiency.
123 Methods like PB2 (Parker-Holder et al., 2020; 2021) use bandit theory to explore hyperparameter
124 spaces, offering performance guarantees, particularly in small population settings. SEARL (Franke
125 et al., 2021) enhances sample efficiency in off-policy RL by employing a shared replay buffer across
126 the population. BG-PBT (Wan et al., 2022) integrates policy distillation (Rusu et al., 2016) to jointly
127 optimize neural architectures and hyperparameters.

128 However, these works do not address a key weakness of PBT: the inherent greediness of its interme-
129 diate selection mechanism. This issue was first identified in the original PBT work (Jaderberg et al.,
130 2017), leading its authors to propose FIRE PBT (Dalibard & Jaderberg, 2021) to mitigate it through
131 learning curve modeling. While FIRE PBT introduces an intricate mechanism (described in section
132 3.1), we aim to introduce a more practical approach of the greediness phenomenon.

133 All the aforementioned approaches rely on a fixed evolution frequency, and do not discuss the choice
134 of the t_{ready} parameter. To our knowledge, we are the first to investigate its impact on PBT, and its
135 connection to greediness.

137 2.2 GREEDINESS AND EVOLUTION FREQUENCY

139 While PBT’s dynamic adaptation of hyperparameters is a key strength, it also introduces a form of
140 greediness in the optimization process. This greediness arises from selecting agents based on their
141 short-term performance, often resulting in an overemphasis on immediate gains at the expense of
142 long-term success. Evolution frequency lies at the core of this problem, as it controls the optimization
143 horizon. Increasing t_{ready} allows PBT to select agents based on longer-term performance, mitigating
144 the short-sighted decisions issue. However, this comes at the cost of sacrificing PBT’s main principle:
145 its dynamic adaptation throughout the training run. We identify two collapse modes that can be
146 caused by too frequent evolution: *diversity collapse* and *hyperparameter collapse*.

147 **Hyperparameter collapse.** Certain hyperparameters, such as the learning rate or exploration
148 factors in RL, are inherently susceptible to greediness. Decaying these hyperparameters often yields
149 immediate performance gains, making them more favorable during short-term selection. However,
150 lower values restrict the exploration of the solution space, reducing the likelihood of finding better
151 optima within t_{ready} steps. This initiates a self-reinforcing cycle: agents with higher learning rates
152 are outperformed and thus replaced by agents with lower learning-rates that fine-tune the found local
153 optimum. After a few evolution steps, this hyperparameter collapse can combine with diversity loss,
154 leading the overall optimization process to a convergence trap.

156 **Diversity collapse.** Diversity loss is a well-known weakness in evolutionary algorithms (EAs,
157 Spears (1995)) that has not been directly addressed in PBT. When optimizing problems with multiple
158 local optima, EAs often lose population diversity and converge to a single basin of attraction.
159 Typically, this issue is corrected using niching techniques (Shir, 2012), which penalize reductions in
160 diversity. In PBT, the repeated cloning of the highest-performing agents at each evolution step leads
161 to a similar problem. Our variance-exploitation experiment in section 4.3 further highlights that this
diversity collapse can cause PBT to fail, independently from hyperparameter optimization.

162 A solution to PBT’s greediness should account for both of these collapses. One straightforward
163 approach is to reduce the evolution frequency, giving agents more time to escape local optima and
164 slowing the loss of diversity. However, this can’t be a satisfactory solution, as it would directly harm
165 PBT’s sample efficiency by allowing poorly performing agents to persist longer. This ultimately
166 pushed PBT closer to a Random Search, where evolution is entirely absent.

168 3 MULTIPLE-FREQUENCIES POPULATION-BASED TRAINING

169
170 To build upon our insights on evolution frequency, we propose MF-PBT, which employs multiple
171 frequencies. By incorporating low-frequency agents that are less susceptible to hyperparameter
172 and diversity collapse, alongside higher-frequency agents that enable quick adaptation and stronger
173 anytime performance, MF-PBT mitigates the greediness of traditional PBT without sacrificing its
174 core strengths.

175 176 3.1 SUB-POPULATIONS

177
178 A key challenge in PBT is the misalignment between short-term and long-term optimization. As the
179 algorithm selects agents solely based on their performance over t_{ready} training steps, and is blind to
180 their long-term potential, it greedily favors hyperparameters that yield immediate gains, eliminating
181 those that could lead to superior performance in the long term. Nevertheless, this short-term feedback
182 is valuable to achieve strong anytime performance.

183 FIRE PBT (Dalibard & Jaderberg, 2021) introduced the concept of using sub-populations to address
184 the trade-off between short-term and long-term optimization. In their approach, one sub-population
185 is allowed to adopt a greedy strategy by directly optimizing the fitness signal, while the others aim
186 to optimize a proxy for long-term performance: the *improvement rate*. To evaluate the long-term
187 potential of hyperparameters, FIRE PBT uses an *evaluator* agent that simulates training with those
188 hyperparameters. The core assumption is that faster improvement in the evaluator’s performance
189 indicates better long-term potential, which is a quite strong assumption on HPO.

190 In contrast, we argue that the best proxy for long-term performance is long-term performance itself.
191 Rather than crafting an estimation, we let some agents train over longer timescales before evolu-
192 tion. In MF-PBT, each sub-population runs PBT at its own distinct evolution frequency. Dynamic
193 sub-populations (i.e., higher frequency) focus on local optimization and short-term improvements,
194 which can be greedy but offer gains in sample efficiency. Conversely, steady (low-frequency) sub-
195 populations assess long-term performance, avoiding the pitfalls of greediness at the expense of sample
196 efficiency.

197 Our main intuition comes from the phenomenon of greediness itself. When an algorithm shows strong
198 early performance but eventually falls behind a simpler baseline, it is a clear sign of over-optimization
199 and entrapment in a poor local solution. Based on this comparison principle, we expect dynamic
200 agents to over-optimize local optima, and use the steady agents to regularly check if the dynamic
201 agents have been greedy. Once greediness is identified, we correct it by restarting the optimization of
202 dynamic agents around a better optimum found by steadier agents, a process managed through our
203 asymmetric migration mechanism, details in next subsection.

204 205 3.2 ASYMMETRIC MIGRATION PROCESS

206 To effectively leverage the sub-populations, instead of running multiple PBT instances independently,
207 an inter-population information transfer mechanism is needed. Alongside the winners, losers, and
208 survivors brackets, MF-PBT introduces a migration bracket, allowing poorly performing agents
209 within a sub-population to be replaced by better-performing agents from other sub-populations. The
210 migration process operates asymmetrically based on the frequencies of the concerned sub-populations.

211 If a dynamic agent is outperformed by an agent from a steadier sub-population, this signals greediness.
212 In response, we replace the dynamic agent with a clone of the steady one, to restore diversity in the
213 dynamic sub-population and avoid convergence traps.

214 Conversely, if a steady agent is outperformed by a more dynamic agent, the dynamic agent’s solution
215 may result from a valuable high-frequency optimization pattern. However, since it might have been

achieved through over-optimization, we protect the steady sub-population from hyperparameter collapse by importing only the dynamic agent’s weights, not its hyperparameters.

3.3 ALGORITHM

Algorithm 1 Multiple-Frequencies Population Based Training (MF-PBT)

```

1: procedure TRAINING( $\mathcal{P}$ )
2:   for  $\delta = 1, \dots, T/t_{\text{ready}}$  do
3:     STEP( $a, \forall a \in \mathcal{P}, t_{\text{ready}}$ ) ▷ Parallel Training for  $t_{\text{ready}}$  steps
4:      $\mathcal{P} \leftarrow \text{RANKING}(\{a \in \mathcal{P}\})$  ▷ Evaluate fitness and sort agents
5:     for  $i = 1, \dots, M$  do
6:       if  $\delta \bmod \delta_i = 0$  then ▷ Population Update
7:          $\mathcal{B}_i \leftarrow \text{BRACKETS}(\mathcal{P}_i)$ 
8:          $\mathcal{P}_i \leftarrow \text{EVOLUTION}(\mathcal{P}_i, \mathcal{B}_i^1, \mathcal{B}_i^4)$ 
9:          $\mathcal{P}_i \leftarrow \text{MIGRATION}(\mathcal{P}_i, \mathcal{P}_{-i}, \mathcal{B}_i^1, \mathcal{B}_i^3)$ 
10:        end if
11:      end for
12:    end for
13: end procedure

```

Similar to PBT, MF-PBT operates with a population of N agents that train concurrently, evaluated every t_{ready} steps and assigned a fitness score. The agents are divided into M sub-populations $\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_M$, each containing $n = N/M$ agents.

Each sub-population \mathcal{P}_i evolves at its own frequency, parameterized by the factor δ_i , meaning it undergoes an evolution step every $\delta_i \times t_{\text{ready}}$ training steps. We set \mathcal{P}_1 to be the reference population, and the δ_i to be integers with $1 = \delta_1 < \delta_2 < \dots < \delta_M$.

Brackets. When a sub-population \mathcal{P}_i evolves, its agents are ranked and divided in four brackets: the top quarter, \mathcal{B}_i^1 (winners); the second quarter, \mathcal{B}_i^2 (survivors); the third quarter, \mathcal{B}_i^3 (open for migration); and the last quarter, \mathcal{B}_i^4 (losers). For simplicity, we assume n is a multiple of 4.

Evolution. Regarding the winners, survivors and losers, MF-PBT behaves identically as PBT. The agents in \mathcal{B}_i^4 (losers) are replaced with perturbed clones of agents from \mathcal{B}_i^1 (winners). The survivors (\mathcal{B}_i^2) continue training unchanged.

Algorithm 2 Asymmetric migration in MF-PBT

```

1: function MIGRATION( $\mathcal{P}_i, \mathcal{P}_{-i}, \mathcal{B}_i^1, \mathcal{B}_i^3$ )
2:    $k = 1$ 
3:   for  $j = 1, \dots, n/4$  do
4:     if  $\text{FITNESS}(\mathcal{B}_i^3(j)) \geq \text{FITNESS}(\mathcal{P}_{-i}(k))$  then
5:       continue ▷ Agents in  $\mathcal{B}_i^3$  better than contenders in  $\mathcal{P}_{-i}$  are kept as is
6:     end if
7:      $i' \leftarrow \text{INDEX}(\mathcal{P}_{-i}(k))$  ▷ Sub-population Index Retrieval
8:     if  $\delta_{i'} < \delta_i$  then
9:        $\mathcal{B}_i^3(j)_\theta \leftarrow \mathcal{P}_{-i}(k)_\theta$  ▷ Weights Assignment
10:       $\mathcal{B}_i^3(j)_h \leftarrow \mathcal{B}_i^1(1)_h$  ▷ Hyperparameter Assignment
11:    else if  $\delta_{i'} > \delta_i$  then
12:       $\mathcal{B}_i^3(j)_{\theta,h} \leftarrow \mathcal{P}_{-i}(k)_{\theta,h}$  ▷ Full Transfer
13:    end if
14:     $k \leftarrow k + 1$ 
15:  end for
16: end function

```

Migration. The agents in \mathcal{B}_i^3 are compared against agents in $\mathcal{P}_{-i} = \mathcal{P} \setminus \mathcal{P}_i$, to determine if they should be replaced by a copy of an external agent. First, both the agents in \mathcal{B}_i^3 and \mathcal{P}_{-i} are sorted in descending order of fitness. Then, we sequentially perform pairwise comparisons of agents in \mathcal{B}_i^3 and

\mathcal{P}_{-i} . For each agent in \mathcal{B}_i^3 , if it is outperformed by the current top external agent, we replace it using the asymmetric logic described in section 3.2. The procedure is detailed in Algorithm 2.

4 EXPERIMENTS

While MF-PBT can be applied to any HPO problem, we focus on reinforcement learning, where its impact is likely most significant. Following Wan et al. (2022), we use the parallelizable Brax framework (Freeman et al., 2021) to train a *Proximal Policy Optimization* (PPO) (Schulman et al., 2017) agent on multiple control tasks.

We use `jax`-based (Bradbury et al., 2018) implementations of MF-PBT and PPO, designed to parallelize agents on GPUs, thereby leveraging the capabilities of the Brax framework. This implementation achieves approximately 10^6 steps per second on two Nvidia A100 40 GB GPUs, allowing us to train over extended timescales and clearly demonstrate PBT’s limitations in the long term. For robust and fair evaluations, we conduct experiments on seven random seeds and report the interquartile means (IQM) (Agarwal et al., 2021) and interquartile ranges (IQR). To ensure reproducibility, we will make our code publicly available.¹

We use a reference value of $t_{\text{ready}} = 10^6$ environment interactions, consistent with BG-PBT’s experiments (Wan et al., 2022). This choice allow us to demonstrate how a conventional value can lead PBT to collapse over extended timescales. For the computation of the fitness score, we evaluate agents on 512 episodes and use the mean evaluation reward. Based on preliminary experiments, we selected $N = 32$ agents split into $M = 4$ sub-populations of $n = 8$ agents each, as moving from 16 to 32 agents significantly improved performance, while gains diminished beyond 32. In this setting, our longest experiments (3 billion steps in the *Humanoid* environment) require approximately 30 hours using two Nvidia A100 40 GB GPUs.

Our computational budget allowed us to train for approximately 1 billion steps per experiment, guiding our choice of $\delta_4 = 50$ for the steadiest sub-population. Indeed, higher values would get it closer to a random search, as the total number of evolution steps for this specific sub-population equals $1000/\delta_4$. To facilitate smoother transfers between the fastest and slowest sub-populations, we selected two intermediary values: $\delta_2 = 10$ and $\delta_3 = 25$. This configuration of the δ -values demonstrated slightly superior performance compared to a less spread geometric progression, as detailed in Appendix B.1. Given that the results already showed MF-PBT’s ability to overcome PBT’s greediness, we did not further tune the δ -values.

We optimize the learning rate and the entropy cost of PPO’s loss (Schulman et al., 2017), as these hyperparameters are particularly susceptible to causing hyperparameter collapse. For all experiments, we initially log-uniformly sample the learning rate between 10^{-5} and 10^{-3} , and the entropy cost between 10^{-3} and 10^{-1} . For the remaining hyperparameters, we use the tuned values proposed by Brax when available.² Notably, the same network architectures are used across all environments.

4.1 COMPARATIVE STUDY OF MF-PBT

We first compare MF-PBT to both PBT and Random Search (RS) (Bergstra & Bengio, 2012), using the same number of agents and the same value of t_{ready} . Since RS does not involve evolution, it can be viewed as a version of PBT with δ set to $+\infty$ (see Appendix A.1 for additional implementation details). This allows us to isolate the effect of evolution in PBT; if RS performs better than PBT, it is a clear sign of greediness.

For the perturbation of hyperparameters in both PBT and MF-PBT, we use the naive *perturb* strategy introduced in the original PBT (Jaderberg et al., 2017), which involves multiplying the hyperparameters by a factor λ randomly sampled from $\{0.8, 1.25\}$. We do not include PB2 (Parker-Holder et al., 2020) and BG-PBT (Wan et al., 2022) as baselines, as these methods build on top of PBT to enhance exploration. Our focus is to identify and mitigate the inherent greediness in PBT’s evolution mechanism. The improvements introduced in MF-PBT could potentially benefit PB2 and BG-PBT as well.

¹The code will be published on GitHub after the double-blind review process. A minimal version of the project is included in the supplementary materials for reviewers.

²See Brax’s GitHub. For *Hopper* and *Walker2D* we used the same values as in *Humanoid*.

FIRE PBT (Dalibard & Jaderberg, 2021) is also not included due to reproducibility challenges. It lacks a public implementation, and key aspects, such as the curve smoothing process, are not detailed in the paper. Additionally, their RL experiments use V-MPO (Song et al., 2019), an algorithm without a public implementation, and their experiment on ImageNet requires 200 TPU-v3 days, making direct comparison prohibitive.

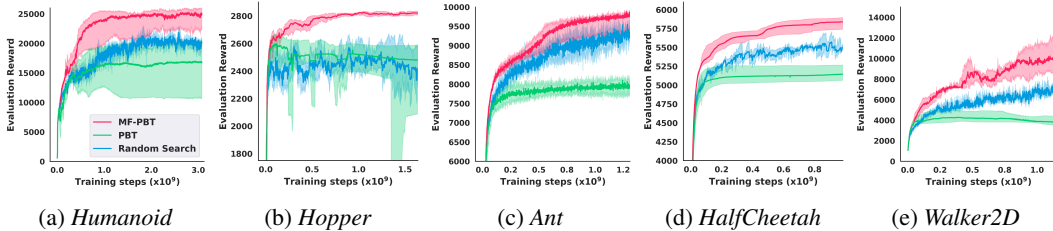


Figure 1: **Performance of MF-PBT, PBT, and RS on Brax environments.** IQM across seven seeds, with IQR shaded. The performance of each algorithm is determined by the highest fitness score (mean evaluation reward over 512 episodes) among the 32 agents, evaluated every t_{ready} training steps.

This experiments highlight the limitations of regular PBT for long-term performance. In Table 1, we report the performance of each algorithm after 50 million training steps, the default horizon proposed by Brax for most environments. At this stage, PBT demonstrates relatively strong performance, achieving results that are superior or comparable to RS in most environments. However, after one billion steps, the same PBT falls significantly behind RS, demonstrating its greediness. Evolving every $t_{\text{ready}} = 10^6$ steps made it collapse in a poor optimum, while RS, which does not evolve, found better solutions.

In contrast, MF-PBT consistently outperforms both PBT and RS at both training horizons, showcasing its adaptability across varying timescales. The training trajectories reported in Figure 1 further illustrate that MF-PBT achieves stronger anytime performance, consistently outperforming RS and PBT throughout the training run. This indicates that MF-PBT has a better sample efficiency, achieving high rewards more rapidly.

Table 1: **IQM of the performance** achieved by the evaluated HPO algorithms at 50 million steps and 1 billion steps across seven random seeds. Methods within the IQR of the best-performing method are bolded. The *PPO* columns correspond to the training of a single agent with the default hyperparameters.

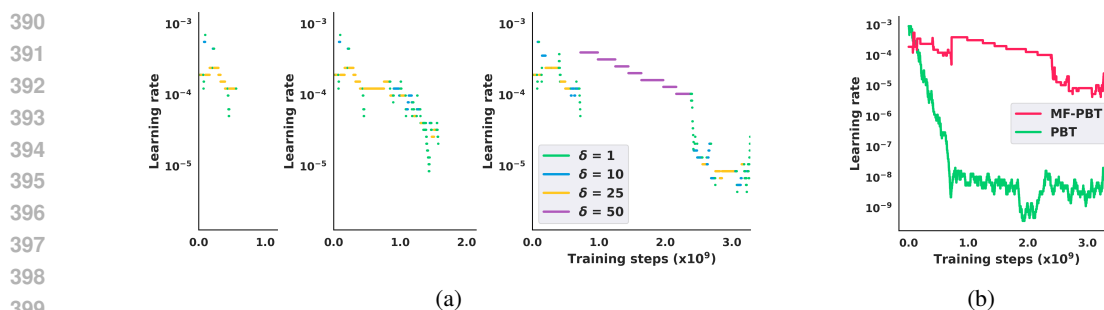
Method	Performance at 50M steps				Performance at 1B steps			
	<i>PPO</i>	RS	PBT	MF-PBT	<i>PPO</i>	RS	PBT	MF-PBT
<i>Humanoid</i>	7903	9021	8348	9266	14934	17713	16171	23793
<i>Hopper</i>	1782	2437	2542	2579	1822	2498	2519	2819
<i>Ant</i>	5482	6858	6820	7115	7102	9050	7900	9654
<i>HalfCheetah</i>	3786	4906	4914	5154	4262	5503	5143	5837
<i>Walker2D</i>	2881	3309	3822	3852	4261	7005	3870	9545

4.2 HYPERPARAMETERS SCHEDULES

To better illustrate MF-PBT’s optimization process, we reconstruct the history of the best agent to visualize its hyperparameter schedule. In figure Figure 2a, we present three snapshots of MF-PBT taken during training on the *Humanoid* environment, at 750 million, 1.5 billion and 3 billion steps. For each snapshot, we trace back the history of the best-performing agent by recursively identifying the agents it cloned. Each colored segment in the schedule indicates the sub-population that produced the agent, showcasing how MF-PBT combines contributions from all sub-populations to produce its final solution.

378 A comparison of the three snapshots shows how MF-PBT is able to target for strong anytime
 379 performance. At every stage of the training, there are greedy agents diving into local optima, in order
 380 to maximize the immediate reward. Steady agents on their side, focus on long-term performance and
 381 protect the overall optimization process from collapse.

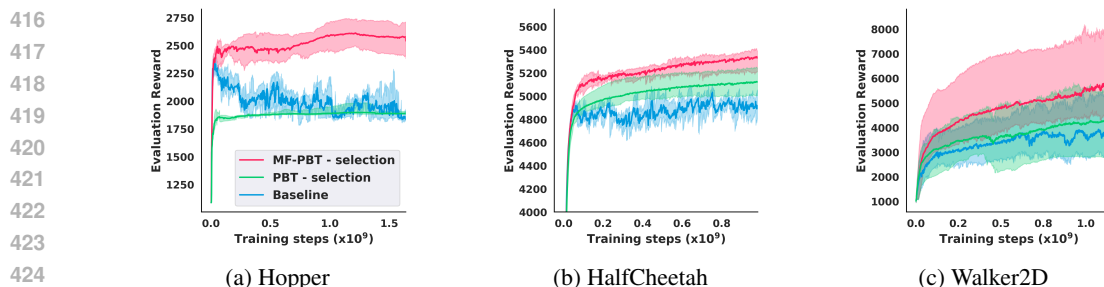
382 In the final schedule, we observe three phases. First, MF-PBT identifies an interesting high-frequency
 383 optimization pattern, where the learning rate increases briefly before decreasing, resembling the
 384 warm-up strategy proposed in Smith (2017). Next, the steady agents, slowly decrease their learning
 385 rate, avoiding collapse and aiming for better long-term rewards. Finally, dynamic agents take the
 386 lead, by fine-tuning the found local optimum through more aggressive learning rate decrease. This
 387 final schedules shows how MF-PBT effectively makes use of its multiple frequencies to produce the
 388 best long-term performance.



400 Figure 2: **Example of learning rate schedules** for MF-PBT and PBT on the *Humanoid* environment.
 401 (a) MF-PBT snapshots at 750 million, 1.5 billion, and 3 billion training steps. Colors represent the
 402 sub-populations contribution to the schedule, showing how MF-PBT integrates input from various
 403 frequencies. (b) Comparison of the two final schedules, illustrating a case of hyperparameter collapse
 404 in PBT.

406 Figure 2b compares the final schedule produced by MF-PBT, to a schedule from a PBT experiment
 407 that encountered a strong hyperparameter collapse, ceasing to improve its reward after only 340
 408 million steps. This collapse results from the presence of strong, peaked local optima in the *Humanoid*
 409 environment, such as running on one leg. Escaping such optima requires extensive exploration, as
 410 deviating from them is highly punitive, leading short-sighted PBT to enter a collapse cycle without
 411 finding better solutions. This difficulty with the *Humanoid* environment has also been noted in
 412 BG-PBT (Wan et al., 2022).

414 4.3 MF-PBT AS A VARIANCE-EXPLOITER



426 Figure 3: **Comparative performance of MF-PBT, PBT and a non-evolutive baseline for variance-**
 427 **exploitation.** IQM across seven seeds, with IQR shaded.

429 Building on our discussion on variance-exploitation in section 2.1, we designed experiments to
 430 evaluate MF-PBT’s ability to leverage stochasticity in training outcomes to improve performance,
 431 even without hyperparameter tuning. In these experiments, all agents are fixed to use the default
 hyperparameters for the entire duration of training, with only weight cloning performed during the

432 evolution steps. To provide a baseline for comparison, we included a non-evolutive approach: running
 433 32 agents independently without any weight replication or hyperparameter tuning, evaluating their
 434 fitness every t_{ready} steps.

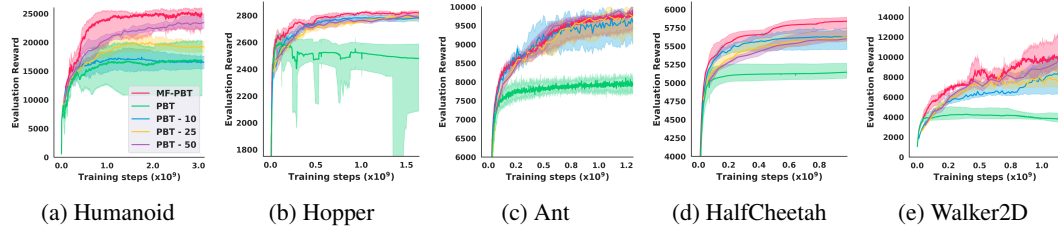
435 The resulting trajectories in Figure 3 reveal three key insights. (1) variance-exploitation can enhance
 436 the performance of a fixed hyperparameter configuration, as demonstrated in the *HalfCheetah*
 437 environment; (2) PBT, even when no hyperparameter collapse is possible, can still fall behind its
 438 non-evolutive counterpart, evidencing diversity collapse- the inherent greediness of the cloning
 439 mechanism; (3) MF-PBT significantly improves performance without modifying hyperparameters,
 440 illustrating the power of a more sophisticated cloning mechanism.

441 Interestingly, while PBT outperformed the non-evolutive baseline in the variance-exploitation regime
 442 for *HalfCheetah* and *Walker2D*, its performance dropped when hyperparameter tuning was introduced,
 443 indicating hyperparameter collapse. In contrast, MF-PBT performed in both regimes, highlighting its
 444 ability to overcome both diversity and hyperparameter collapse.

447 5 ABLATIVE STUDIES

449 5.1 EVOLUTION FREQUENCY

451 Our intuition is that evolving less frequently (increasing δ) mitigates greediness and ensures better
 452 long-term performance, but using multiple frequencies is necessary to achieve stronger anytime
 453 performance. To test this, we conducted an experiment comparing MF-PBT with four separate
 454 PBT runs, each using 32 agents and evolving at one of the frequencies used within MF-PBT:
 455 $\delta \in \{1, 10, 25, 50\}$.



465 **Figure 4: Impact of the evolution frequency in PBT.** IQM across seven seeds, with IQR shaded.

468 The resulting trajectories plotted in Figure 4 confirm our first intuition about the critical role of
 469 evolution frequency, demonstrating its significant impact on PBT’s performance. The curves also
 470 reveal that the best frequencies vary by task; for example, on *Humanoid*, $\delta = 50$ is the most
 471 effective, whereas on *HalfCheetah*, $\delta = 25$ yields better results. Additionally, most of the slower
 472 PBT configurations outperform RS, indicating that $\delta = +\infty$ is sub-optimal. This underscores the
 473 brittleness of population-based approaches to the choice of t_{ready} .

474 In contrast, MF-PBT achieves either superior or comparable final performance relative to each
 475 single-frequency PBT experiment, while also offering significant sample efficiency gains in most
 476 environments. This indicates that employing multiple frequencies within MF-PBT is superior to
 477 relying on a single frequency. Moreover, MF-PBT’s ability to outperform each of its sub-components
 478 simplifies the selection of δ -values, as MF-PBT will always perform at least as well as its best-
 479 performing sub-population.

481 5.2 SYMMETRIC MIGRATION

483 We now assess the importance of the asymmetry in the migration process, which adds a protection
 484 against hyperparameter collapse by preventing greedy agents from corrupting steadier sub-populations.
 485 To test this, we compare MF-PBT with an alternative version where hyperparameters are always
 transferred along with weights, regardless of the δ -values.

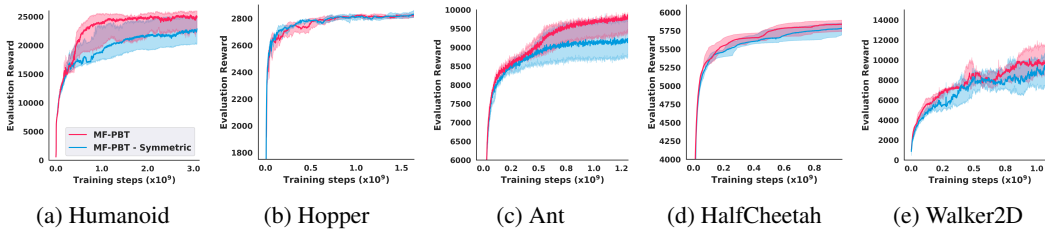


Figure 5: Ablation on the asymmetric migration. IQM across seven seeds, with IQR shaded.

The training trajectories in Figure 5 show that while the asymmetry has little impact on *Hopper*, it yields improvements in most environments, particularly in the challenging *Humanoid* task. This indicates that the asymmetric design indeed enhances long-term performance.

6 CONCLUSION

We introduced MF-PBT, an extension of Population-Based Training, designed to address the inherent greediness in traditional PBT. Our experiments on various reinforcement learning tasks identified two key failure modes of PBT: diversity and hyperparameter collapse, both linked to the evolution frequency. Building on these insights, we proposed MF-PBT, which incorporates multiple sub-populations evolving at different frequencies and an asymmetric migration process to balance short and long-term optimization. The results demonstrated that MF-PBT effectively overcomes both collapses associated with PBT while maintaining strong anytime performance.

Through ablation studies, we highlighted the critical role of evolution frequency in PBT and showed that using multiple frequencies increases robustness to this parameter. We believe this insight could be broadly valuable for all population-based approaches. Combining MF-PBT’s mechanisms with other extensions to PBT, such as PB2 (Parker-Holder et al., 2020), PB2-Mix (Parker-Holder et al., 2021), or BG-PBT (Wan et al., 2022), could further enhance performance. Further work could also include a study of MF-PBT’s performance on other tasks such as supervised learning.

Our experiments about variance-exploitation highlight that a non-negligible share of the performance gains in population-based methods arises from leveraging exploration luck rather than tuning hyperparameters effectively. This underscores the need for a more comprehensive study on the origins of improvements brought by population-based methods.

REFERENCES

- Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Belle-mare. Deep reinforcement learning at the edge of the statistical precipice. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 29304–29320. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/f514cec81cb148559cf475e7426eed5e-Paper.pdf.
- Noor Awad, Neeratoy Mallik, and Frank Hutter. Dehb: Evolutionary hyberband for scalable, robust and efficient hyperparameter optimization. In Zhi-Hua Zhou (ed.), *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pp. 2147–2153. International Joint Conferences on Artificial Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/296. URL <https://doi.org/10.24963/ijcai.2021/296>. Main Track.
- Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi, Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the Atari human benchmark. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 507–517. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/badia20a.html>.

-
- 540 James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal*
541 *of Machine Learning Research*, 13(10):281–305, 2012. URL [http://jmlr.org/papers/](http://jmlr.org/papers/v13/bergstra12a.html)
542 [v13/bergstra12a.html](http://jmlr.org/papers/v13/bergstra12a.html).
543
- 544 James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-
545 parameter optimization. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Wein-
546 berger (eds.), *Advances in Neural Information Processing Systems*, volume 24. Curran Asso-
547 ciates, Inc., 2011. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf)
548 [2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf).
- 549 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
550 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and
551 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL
552 <http://github.com/google/jax>.
553
- 554 Valentin Dalibard and Max Jaderberg. Faster improvement rate population based training. *CoRR*,
555 [abs/2109.13800](https://arxiv.org/abs/2109.13800), 2021. URL <https://arxiv.org/abs/2109.13800>.
- 556 Theresa Eimer, Marius Lindauer, and Roberta Raileanu. Hyperparameters in reinforcement learning
557 and how to tune them. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt,
558 Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on*
559 *Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 9104–9149.
560 PMLR, 23–29 Jul 2023. URL [https://proceedings.mlr.press/v202/eimer23a.](https://proceedings.mlr.press/v202/eimer23a.html)
561 [html](https://proceedings.mlr.press/v202/eimer23a.html).
562
- 563 Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: Robust and efficient hyperparameter opti-
564 mization at scale. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International*
565 *Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*,
566 pp. 1437–1446. PMLR, 10–15 Jul 2018. URL [https://proceedings.mlr.press/v80/](https://proceedings.mlr.press/v80/falkner18a.html)
567 [falkner18a.html](https://proceedings.mlr.press/v80/falkner18a.html).
568
- 569 Matthias Feurer and Frank Hutter. Hyperparameter optimization. In *Automatic Machine Learning:*
570 *Methods, Systems, Challenges*, pp. 3–38. Springer, 2019.
- 571 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of
572 deep networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International*
573 *Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp.
574 1126–1135. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.press/v70/](https://proceedings.mlr.press/v70/finn17a.html)
575 [finn17a.html](https://proceedings.mlr.press/v70/finn17a.html).
576
- 577 Jörg K.H. Franke, Gregor Koehler, André Biedenkapp, and Frank Hutter. Sample-efficient automated
578 deep reinforcement learning. In *International Conference on Learning Representations*, 2021.
579 URL <https://openreview.net/forum?id=hSjxQ3B7GWq>.
580
- 581 C. Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem.
582 Brax – a differentiable physics engine for large scale rigid body simulation, 2021.
- 583 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.
584 Deep reinforcement learning that matters. In *Proceedings of the Thirty-Second AAAI Confer-*
585 *ence on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence*,
586 *AAAI’18/IAAI’18/EAAI’18*. AAAI Press, 2018. ISBN 978-1-57735-800-8.
587
- 588 Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M. Czarnecki, Jeff Donahue, Ali
589 Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, Chrisantha Fernando, and
590 Koray Kavukcuoglu. Population based training of neural networks. *CoRR*, [abs/1711.09846](https://arxiv.org/abs/1711.09846), 2017.
591 URL <http://arxiv.org/abs/1711.09846>.
592
- 593 Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband:
A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning*
Research, 18(185):1–52, 2018. URL [http://jmlr.org/papers/v18/16-558.html](http://jmlr.org/papers/v18/li18a.html).

594 Siqi Liu, Guy Lever, Zhe Wang, Josh Merel, S. M. Ali Eslami, Daniel Hennes, Wojciech M. Czarnecki,
595 Yuval Tassa, Shayegan Omidshafiei, Abbas Abdolmaleki, Noah Y. Siegel, Leonard Hasenclever,
596 Luke Marris, Saran Tunyasuvunakool, H. Francis Song, Markus Wulfmeier, Paul Muller, Tuomas
597 Haarnoja, Brendan Tracey, Karl Tuyls, Thore Graepel, and Nicolas Heess. From motor control
598 to team play in simulated humanoid football. *Science Robotics*, 7(69):eabo0235, 2022. doi:
599 10.1126/scirobotics.abo0235. URL [https://www.science.org/doi/abs/10.1126/
600 scirobotics.abo0235](https://www.science.org/doi/abs/10.1126/scirobotics.abo0235).

601 Jack Parker-Holder, Vu Nguyen, and Stephen J Roberts. Provably efficient online hyperparameter opti-
602 mization with population-based bandits. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and
603 H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 17200–17211.
604 Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_
605 files/paper/2020/file/c7af0926b294e47e52e46cfebe173f20-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/c7af0926b294e47e52e46cfebe173f20-Paper.pdf).

606
607 Jack Parker-Holder, Vu Nguyen, Shaan Desai, and Stephen J Roberts. Tuning mixed in-
608 put hyperparameters on the fly for efficient population based autorl. In M. Ranzato,
609 A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neu-
610 ral Information Processing Systems*, volume 34, pp. 15513–15528. Curran Associates, Inc.,
611 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/
612 file/82debd8a12b498e765a11a8e51159440-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/82debd8a12b498e765a11a8e51159440-Paper.pdf).

613 Jack Parker-Holder, Raghu Rajan, Xingyou Song, André Biedenkapp, Yingjie Miao, Theresa Eimer,
614 Baohe Zhang, Vu Nguyen, Roberto Calandra, Aleksandra Faust, Frank Hutter, and Marius Lindauer.
615 Automated reinforcement learning (autorl): A survey and open problems. *Journal of Artificial
616 Intelligence Research*, 74:517–568, 06 2022. doi: 10.1613/jair.1.13596.

617 Andrei A. Rusu, Sergio Gomez Colmenarejo, Çağlar Gülçehre, Guillaume Desjardins, James Kirk-
618 patrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. Policy
619 distillation. In Yoshua Bengio and Yann LeCun (eds.), *4th International Conference on Learning
620 Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*,
621 2016. URL <http://arxiv.org/abs/1511.06295>.

622 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
623 optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL [http://dblp.uni-trier.
624 de/db/journals/corr/corr1707.html#SchulmanWDRK17](http://dblp.uni-trier.de/db/journals/corr/corr1707.html#SchulmanWDRK17).

625
626 Ofer M. Shir. *Niching in Evolutionary Algorithms*, pp. 1035–1069. Springer Berlin Heidelberg,
627 Berlin, Heidelberg, 2012. ISBN 978-3-540-92910-9. doi: 10.1007/978-3-540-92910-9_32. URL
628 https://doi.org/10.1007/978-3-540-92910-9_32.

629 Leslie N. Smith. Cyclical learning rates for training neural networks. In *2017 IEEE Winter Conference
630 on Applications of Computer Vision (WACV)*, pp. 464–472, 2017. doi: 10.1109/WACV.2017.58.

631
632 H. Francis Song, Abbas Abdolmaleki, Jost Tobias Springenberg, Aidan Clark, Hubert Soyer,
633 Jack W. Rae, Seb Noury, Arun Ahuja, Siqi Liu, Dhruva Tirumala, Nicolas Heess, Dan Belov,
634 Martin A. Riedmiller, and Matthew M. Botvinick. V-MPO: on-policy maximum a posteriori
635 policy optimization for discrete and continuous control. *CoRR*, abs/1909.12238, 2019. URL
636 <http://arxiv.org/abs/1909.12238>.

637 William Michael Spears. Adapting Crossover in Evolutionary Algorithms. In *Evolutionary Pro-
638 gramming IV: Proceedings of the Fourth Annual Conference on Evolutionary Programming*.
639 The MIT Press, 08 1995. ISBN 9780262290920. doi: 10.7551/mitpress/2887.003.0035. URL
640 <https://doi.org/10.7551/mitpress/2887.003.0035>.

641
642 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press,
643 second edition, 2018. URL [http://incompleteideas.net/book/the-book-2nd.
644 html](http://incompleteideas.net/book/the-book-2nd.html).

645 Xingchen Wan, Cong Lu, Jack Parker-Holder, Philip J. Ball, Vu Nguyen, Binxin Ru, and
646 Michael Osborne. Bayesian generational population-based training. In Isabelle Guyon, Mar-
647 rius Lindauer, Mihaela van der Schaar, Frank Hutter, and Roman Garnett (eds.), *Proceed-
ings of the First International Conference on Automated Machine Learning*, volume 188 of

648 *Proceedings of Machine Learning Research*, pp. 14/1–27. PMLR, 25–27 Jul 2022. URL
649 <https://proceedings.mlr.press/v188/wan22a.html>.
650

651 Zhongwen Xu, Hado P van Hasselt, and David Silver. Meta-gradient reinforcement learn-
652 ing. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Gar-
653 nett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Asso-
654 ciates, Inc., 2018. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2018/file/2715518c875999308842e3455eda2fe3-Paper.pdf)
655 [2018/file/2715518c875999308842e3455eda2fe3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/2715518c875999308842e3455eda2fe3-Paper.pdf).

656 Baohe Zhang, Raghu Rajan, Luis Pineda, Nathan Lambert, André Biedenkapp, Kurtland Chua, Frank
657 Hutter, and Roberto Calandra. On the importance of hyperparameter optimization for model-
658 based reinforcement learning. In Arindam Banerjee and Kenji Fukumizu (eds.), *Proceedings*
659 *of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of
660 *Proceedings of Machine Learning Research*, pp. 4015–4023. PMLR, 13–15 Apr 2021. URL
661 <https://proceedings.mlr.press/v130/zhang21n.html>.
662

663 A ADDITIONAL IMPLEMENTATION DETAILS

664 A.1 RANDOM SEARCH BASELINE

665
666 In our experiments, Random Search (RS) serves as a simple baseline for hyperparameter optimization.
667 RS involves randomly sampling hyperparameter values at the start of training and keeping these
668 values fixed throughout the entire training process. Unlike PBT, RS does not involve any evolution
669 or adjustment of hyperparameters based on intermediate performance. Instead, the goal of RS
670 is to evaluate different fixed hyperparameter configurations by following their reward curves and
671 identifying which sampled configuration performs best.
672

673 For this comparison, hyperparameters in RS were sampled uniformly from the same search space as
674 PBT and MF-PBT. By comparing RS to PBT, we isolate the impact of PBT’s evolutionary process;
675 if RS outperforms PBT, it indicates that evolving too frequently can lead to suboptimal long-term
676 performance, which we refer to as the greediness issue.
677

678 A.2 PBT’S PARAMETERS

679
680 In subsection 2.2 we identified that the main source of the greediness issue is that agents do not
681 survive long enough to escape poor local optima and maintain diversity. Alongside the evolution
682 frequency, another parameter of PBT impact the lifespan of agents in the population: the selection
683 rate in the exploit phase.

684 Indeed, in PBT, at each evolution step, 25% of the population is discarded and replaced by copied of
685 the top-agents. One could play on this parameter to mitigate greediness, and create a method similar
686 to MF-PBT, where each sub-population would have its own selection rate. However, as we identify
687 the issue to be about the lifespan of agents, and optimizing for various horizons, we found more
688 natural to frame it explicitly in terms of evolution frequency.

689 We decided to use standard values for the *exploit* and *explore* process of PBT, and keep the same
690 values for MF-PBT in order to isolate the impact of evolution frequency.
691

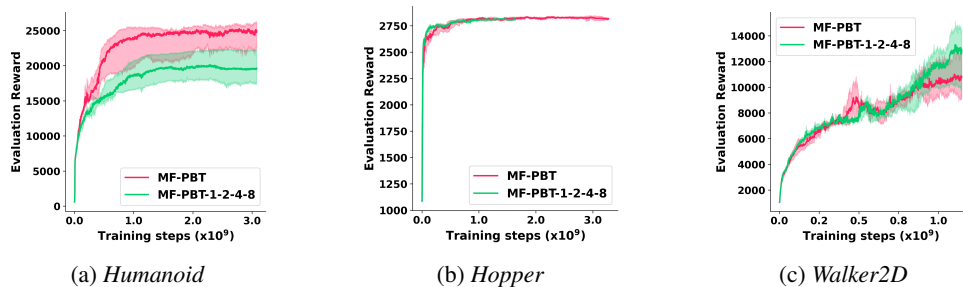
692 B CHOICE OF PARAMETERS

693 B.1 FREQUENCIES

694
695 Figure 6 compares our chosen configuration ($t_{\text{ready}} = 10^6$, $\delta_1 = 1$, $\delta_2 = 10$, $\delta_3 = 25$, $\delta_4 = 50$) with
696 an alternative setup using a geometric progression ($t_{\text{ready}} = 6 \times 10^6$, $\delta_1 = 1$, $\delta_2 = 2$, $\delta_3 = 4$, $\delta_4 = 8$).
697 The goal of this comparison is to assess how the spread of δ -values impacts MF-PBT’s performance.
698

699 While the geometric progression shows a slight advantage on *Hopper* and *Walker2D*, it performs
700 significantly worse on *Humanoid*. Therefore, we opted to continue using the more spread-out
701 configuration.

702
703
704
705
706
707
708
709
710
711



712 **Figure 6: Comparative performance of two configurations for MF-PBT.** IQM across seven seeds,
713 with IQR shaded.

714
715
716

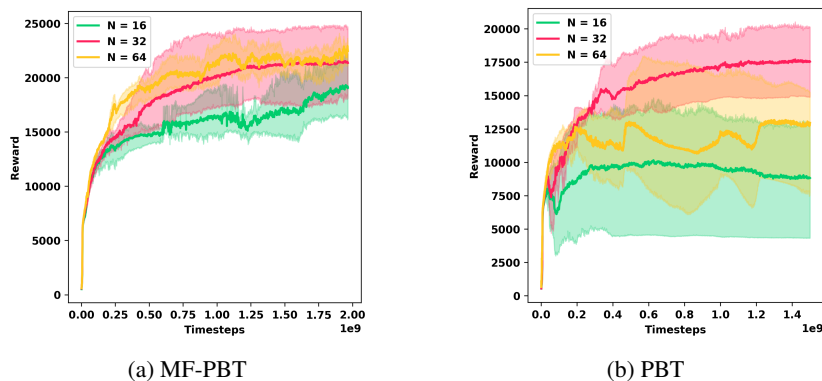
B.2 POPULATION SIZE

717
718
719
720
721
722
723
724

To make a choice for N after fixing the δ -values, we conducted a preliminary experiment on *Humanoid*, the most computationally demanding environment. As shown in Figure 7, the gain from rising from $N = 16$ to $N = 32$ is quite large for both methods. While increasing from $N = 32$ to $N = 64$ was still beneficial for MF-PBT, but the with a much smaller gap.

Interestingly, PBT’s performance decreases with 64 agents on the *Humanoid* task, likely due to the abundance of local optima. With a large population, PBT may quickly converge on a high-performing local optimum, which then limits further exploration.

725
726
727
728
729
730
731
732
733
734
735
736
737



738 **Figure 7: Impact of the population size.** IQM across five seeds, with IQR shaded. Experiments on
739 the *Humanoid* environment.

740
741
742

C ADDITIONNAL EXPERIMENTS

743
744

C.1 PUSHER ENVIRONMENT

745
746
747
748

We made an experiment in the *Pusher* environment from Brax, keeping the same parameters for MF-PBT, PBT and RS and report the training curves in Figure 8

C.2 INCREASING POPULATION SIZE

749
750
751
752
753
754
755

One solution to improve PBT’s performance can be to increase the population size. In Jaderberg et al. (2017), a value of $N = 80$ was used. To make sure we didn’t unfairly treat PBT by picking $N = 32$, we made an additional experiment to compare the gains of using MF-PBT to the gains of simply increasing N in standard PBT.

The curves in Figure 9 show that while on *Hopper* raising to 80 agents greatly improves PBT’s performance, it is not sufficient in a more complex locomotion environment like *Walker2D*.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

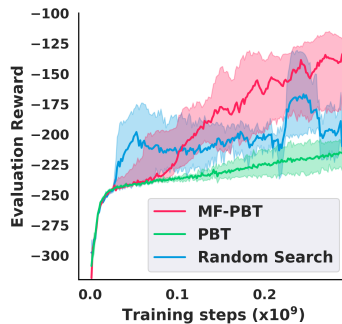


Figure 8: Performance of MF-PBT, PBT, and RS on Pusher. IQM across seven seeds, with IQR shaded.

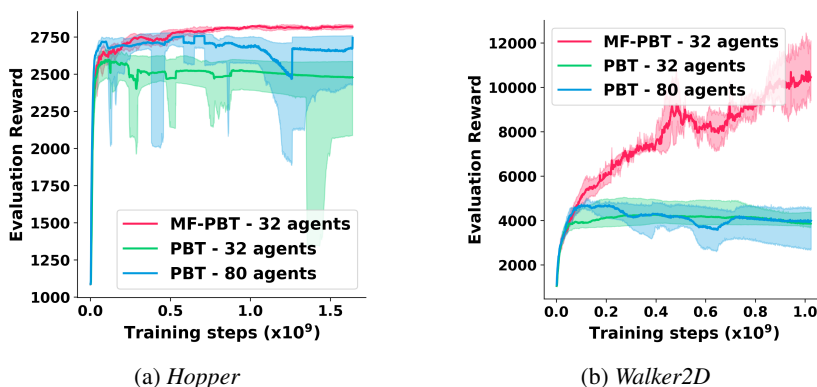


Figure 9: Increasing population size. IQM across seven seeds, with IQR shaded.

C.3 BACKTRACKING

Zhang et al. (2021) proposed to add a backtracking mechanism to PBT, to prevent it from catastrophic forgetting. The method, dubbed PBT-BT (PBT with backtracking), keeps track of the N_e best agents encountered during the training: the *elites*. And every δ evolution steps, the elites are reincorporated into the population.

Since in the *Hopper* and *Humanoid* environments, we observed a substantial amount of runs where PBT’s performance would dramatically drop, PBT-BT could be an interesting alternative baseline in those environments.

The backtracking can be seen as a migration across times, where elites from the past are reincorporated in the population, to enable it to resume training from a better checkpoint. However there is one fundamental difference, in PBT-BT the elites come from the past and didn’t interact as much with the environment; whereas in MF-PBT the steady agents that migrate are "current" agents, meaning they performed the exact same amount of training steps. In MF-PBT, the agents that migrate only differ on their HPO-objective, e.g. performance on 50M steps instead of performance on 1M steps. While backtracking enables recovering from collapses, there is no notion of increasing the lifespan of some hyperparameters to assess their long-term performance.

We implemented PBT-BT with $N = 32$, $N_e = 16$ and $\delta = 50$. The training curves in Figure 10 shows that it improves PBT on *Hopper* by correcting the catastrophic forgetting behavior. However on *Humanoid*, the elites tend to rapidly all belong to the same local optimum, and then PBT-BT is stuck without being able to explore for better solution.

In both cases, MF-PBT outperforms PBT-BT, highlighting that backtracking is not sufficient to overcome PBT’s greediness.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

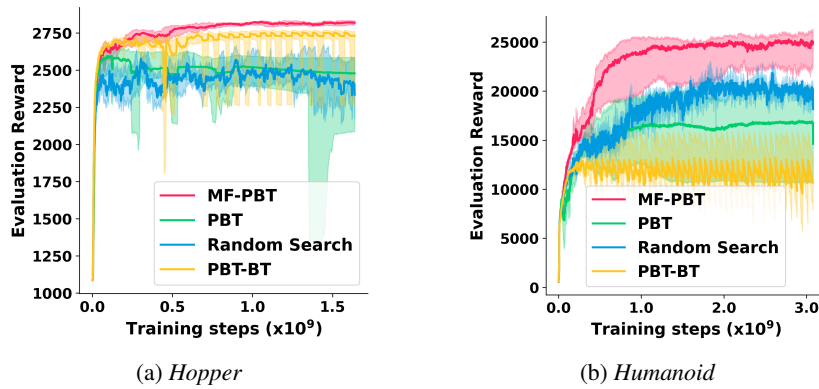


Figure 10: Comparative performance of PBT-BT. IQM across seven seeds, with IQR shaded.