

---

# Diversity Policy Gradient for Sample Efficient Quality-Diversity Optimization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 A fascinating aspect of nature lies in its ability to produce a large and diverse  
2 collection of organisms that are all high-performing in their niche. By contrast,  
3 most AI algorithms focus on finding a single efficient solution to a given problem.  
4 Aiming for diversity in addition to performance is a convenient way to deal with the  
5 exploration-exploitation trade-off that plays a central role in learning. It also allows  
6 for increased robustness when the returned collection contains several working  
7 solutions to the considered problem, making it well-suited for real applications such  
8 as robotics. Quality-Diversity (QD) methods are evolutionary algorithms designed  
9 for this purpose. This paper proposes a novel algorithm, QD-PG, which combines  
10 the strength of Policy Gradient algorithms and Quality Diversity approaches to  
11 produce a collection of diverse and high-performing neural policies in continuous  
12 control environments. The main contribution of this work is the introduction of a  
13 Diversity Policy Gradient (DPG) that exploits information at the time-step level to  
14 thrive policies towards more diversity in a sample-efficient manner. Specifically,  
15 QD-PG selects neural controllers from a MAP-Elites grid and uses two gradient-  
16 based mutation operators to improve both quality and diversity, resulting in stable  
17 population updates. Our results demonstrate that QD-PG generates collections of di-  
18 verse solutions that solve challenging exploration and control problems while being  
19 two orders of magnitude more sample-efficient than its evolutionary competitors.

## 20 1 Introduction

21 Natural evolution has the fascinating ability to produce diverse organisms that are all well adapted to  
22 their respective niche. Inspired by this ability to produce a tremendous diversity of living systems,  
23 Quality-Diversity (QD) is a new family of optimization algorithms that aims at searching for a  
24 collection of both diverse and high-performing solutions (Pugh et al., 2016; Cully & Demiris, 2017).  
25 While classic optimization methods focus on finding a single efficient solution, QD optimization aims  
26 to cover the range of possible solution types and to return the best solution for each type. This process  
27 is sometimes referred to as “illumination” in opposition to optimization, as it reveals (or illuminates)  
28 a search space of interest often called the *behavior descriptor space* (Mouret & Clune, 2015).

29 The principal advantage of QD approaches resides in their intrinsic capacity to deliver a large and  
30 diverse set of working alternatives when a single solution fails (Cully et al., 2015). By producing a  
31 collection of solutions instead of a unique one, QD algorithms allow to obtain different ways to solve  
32 a single problem, leading to greater robustness, which can help to reduce the reality gap when applied  
33 to robotics (Koos et al., 2012). Diversity seeking is the core component that allows QD algorithms to  
34 generate large collections of diverse solutions. By encouraging the emergence of novel behaviors in  
35 the population without focusing on performance alone, diversity seeking algorithms explore regions  
36 of the behavior descriptor space that are unreachable for conventional algorithms (Doncieux et al.,



Figure 1: The agent robot is rewarded for running forward as fast as possible. Following the reward signal without further exploration leads the agent into the trap, which corresponds to a poor local minimum. QD-PG produces a collection of solutions that are diverse and high-performing, allowing to find several working alternatives to solve a deceptive control problem.

37 2019). Another benefit of QD is its ability to solve hard exploration problems where the reward signal  
 38 is sparse or deceptive, and on which standard optimization techniques are ineffective (Colas et al.,  
 39 2020). This ability can be interpreted as a direct consequence of the structured search for diversity in  
 40 the behavior descriptor space.

41 Quality-Diversity algorithms build on black-box optimization methods such as evolutionary algo-  
 42 rithms to evolve a population of solutions (Cully & Demiris, 2017). Historically, they rely on random  
 43 mutations to explore small search spaces but struggle when facing higher-dimensional problems. As  
 44 a result, they often scale poorly to problems where neural networks with many parameters provide  
 45 state-of-the-art results (Colas et al., 2020).

46 Building large and efficient controllers that work with continuous actions has been a long-standing  
 47 goal in Artificial Intelligence and in particular in robotics. Deep reinforcement learning (RL), and  
 48 especially Policy Gradient (PG) methods have proven efficient at training such large controllers  
 49 (Schulman et al., 2017; Lillicrap et al., 2015; Fujimoto et al., 2018; Haarnoja et al., 2018). One of the  
 50 keys to this success lies in the fact that PG methods exploit the structure of the objective function  
 51 when the problem can be formalized as a Markov Decision Process (MDP), leading to substantial  
 52 gains in sample efficiency. Moreover, they also exploit the analytical structure of the controller when  
 53 known, which allows the sample complexity of these methods to be independent of parameter space  
 54 dimensionality (Vemula et al., 2019). In real-world applications, these gains turn out to be critical  
 55 when interacting with the environment is expensive. PG methods usually rely on simple exploration  
 56 mechanisms, like adding Gaussian noise (Fujimoto et al., 2018) or maximizing entropy (Haarnoja  
 57 et al., 2018) to explore the action space, which happens to be insufficient in hard exploration tasks  
 58 where the reward signal is sparse or deceptive (Colas et al., 2018; Nasiriany et al., 2019).

59 Successful attempts have been made to combine evolutionary methods and reinforcement learning  
 60 (Khadka et al., 2019; Khadka & Tumer, 2018; Pourchot & Sigaud, 2018; Shi et al., 2020). However,  
 61 all these techniques only focus on building high-performing solutions and do not explicitly encourage  
 62 diversity within the population. In this regard, they fail when confronted with hard exploration  
 63 problems. To address these problems, one needs to seek both high-performing solutions and diversity  
 64 within them.

## 65 Contributions

66 In this work, we introduce the idea of a *diversity policy gradient* (DPG) that thrives solutions towards  
 67 more diversity. We show that the DPG can be used in combination with the standard policy gradient,  
 68 dubbed *quality policy gradient* (QPG), to produce high-performing and diverse solutions. Our  
 69 algorithm, called QD-PG, builds on MAP-Elites (Mouret & Clune, 2015), demonstrates remarkable  
 70 sample efficiency brought by off-policy PG methods, and produces collections of good solutions  
 71 in a single run (see Figure 1). We compare QD-PG to state-of-the-art RL algorithms and to several  
 72 evolutionary methods known as Evolution Strategies (ESs) augmented with a diversity objective,

73 namely the NS-ES family (Conti et al., 2018) and the ME-ES algorithm (Colas et al., 2020). We  
 74 show that QD-PG generates collections of robust solutions in hard exploration problems while RL  
 75 algorithms struggle to produce a single one, and that QD-PG is two orders of magnitude more sample  
 76 efficient than the best of its evolutionary competitors.

## 77 2 Background

### 78 Problem statement

79 We consider an MDP  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma)$  where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  the action space,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$   
 80 the reward function,  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  the dynamics transition function and  $\gamma$  a discount factor.  
 81 We assume that both  $\mathcal{S}$  and  $\mathcal{A}$  are continuous and consider a controller, or policy,  $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$   
 82 parameterized by  $\theta \in \Theta$ , which is called a *solution* to the problem. We say that a solution  $\theta$  is  
 83 *highly-performing* if the expectation over the sum of rewards is high when using  $\pi_\theta$ . The *fitness* of a  
 84 solution measures its performance  $F : \Theta \rightarrow \mathbb{R}$  where  $F(\theta) = \mathbb{E}_{\pi_\theta} \sum_t \gamma^t r_t$ .

85 To characterize the novelty of a solution w.r.t.  $J$  other solutions, as in QD methods, we introduce a  
 86 behavior descriptor (BD) space  $\mathcal{B}$ , a behavior descriptor extraction function  $\xi : \Theta \rightarrow \mathcal{B}$ , and define a  
 87 distance metric  $\|\cdot\|_{\mathcal{B}}$  over  $\mathcal{B}$ . The *novelty*  $n : \Theta \times \Theta^J \rightarrow \mathbb{R}^+$  of a solution  $\theta$  w.r.t. a list of solutions  
 88  $(\theta_j)_{j=1, \dots, J}$  is defined as  $n(\theta, (\theta_j)_{j=1, \dots, J}) = \sum_j \|\xi(\theta), \xi(\theta_j)\|_{\mathcal{B}}$ . In other words, we quantify the  
 89 novelty of a solution w.r.t. a list of  $J$  solutions as the sum of distances between its behavior descriptor  
 90 and the behavior descriptors of all solutions of the list. We also use the distance  $\|\cdot\|_{\mathcal{B}}$  to characterize  
 91 the *diversity* of a set of  $K$  solutions  $\{\theta_k\}_{k=1, \dots, K}$ . We formally define diversity  $d : \Theta^K \rightarrow \mathbb{R}^+$  as

$$d(\{\theta_k\}_{k=1, \dots, K}) = \sum_{i=1}^K \min_{k \neq i} \|\xi(\theta_i), \xi(\theta_k)\|_{\mathcal{B}}, \quad (1)$$

92 meaning that a set of solutions is diverse if the solutions are distant with respect to each other in the  
 93 sense of  $\|\cdot\|_{\mathcal{B}}$ .

### 94 The MAP-Elites algorithm

95 MAP-Elites (Mouret & Clune, 2015) is a simple yet state-of-the-art QD algorithm that has been  
 96 successfully applied to a wide range of challenging problems such as robot damage recovery (Cully  
 97 et al., 2015), molecular robotic control (Cazenille et al., 2019) and game design (Alvarez et al., 2019).  
 98 In MAP-Elites, the behavior descriptor space  $\mathcal{B}$  is discretized into a grid of cells, also called niches,  
 99 with the aim of filling each cell with a high-performing solution. The algorithm starts with an empty  
 100 grid and an initial random set of  $K$  solutions that are evaluated and added to the grid by following  
 101 simple insertion rules. If the cell corresponding to the behavior descriptors of a solution is empty, then  
 102 the solution is added to this cell. If there is already a solution in the cell, the new solution replaces it  
 103 only if it has greater fitness. At each iteration,  $P$  existing solutions are sampled uniformly from the  
 104 grid and randomly mutated to create  $P$  new solutions. These new solutions are then evaluated and  
 105 added to the grid following the same insertion rules. This cycle is repeated until convergence or for a  
 106 given budget of iterations.

107 Though MAP-Elites is a compelling and efficient method, it suffers from a low sample efficiency  
 108 as it relies on random mutations. Recently, Colas et al. (2020) tackled this problem by updating  
 109 the solutions through an Evolution Strategy known as the Cross-Entropy method. Notably, they  
 110 showed that MAP-Elites could be scaled with their method to address complex MUJOCO control  
 111 environments at the cost of very large computational resources. In this study, we propose to harness  
 112 policy gradients (QPG and DPG) to build a more sample-efficient MAP-Elites approach.

## 113 3 Key Principle: Diversity Policy Gradient

114 Let us assume that we have a MAP-Elites grid containing  $K$  solutions  $(\theta_1, \dots, \theta_K)$ . To increase  
 115 diversity in the grid using the DPG, we need to update one sampled solution  $\theta$  from the grid using  
 116 gradient ascent. To do so, we aim to compute the gradient of the population diversity w.r.t.  $\theta$ , where  
 117 diversity is defined in Equation (1). As the  $K$  solutions are independent, order does not matter and

118 we can consider optimizing arbitrarily  $\theta = \theta_1$ . To compute the gradient of  $d$  w.r.t.  $\theta_1$ , we need  
 119 to separate the terms that depend on  $\theta_1$  from the others. The terms that depend on  $\theta_1$  correspond  
 120 to the distance of  $\theta_1$  to its nearest neighbor, which we define as  $\theta_2$ , and to the distances of  $\theta_1$  to  
 121 the  $\theta$ s for which  $\theta_1$  is the nearest neighbor. We can arbitrarily index them from 3 to  $J^1$ , thus:

$$122 \quad d(\{\theta_k\}_{k=1,\dots,K}) = \sum_{j=2}^J \|\xi(\theta_1), \xi(\theta_j)\|_{\mathcal{B}} + M, \text{ where } M = \sum_{i \notin \{1,\dots,J\}} \min_{k \neq i} \|\xi(\theta_i), \xi(\theta_k)\|_{\mathcal{B}}.$$

123 Only the first term of the sum depends on  $\theta = \theta_1$ . Furthermore, we observe that this term equals the  
 124 novelty of solution  $\theta_1$  w.r.t. the list  $(\theta_j)_{2 \leq j \leq J}$ . Therefore, the gradient of diversity w.r.t.  $\theta_1$  is

125  $\nabla_{\theta_1} d(\{\theta_k\}_{k=1,\dots,K}) = \nabla_{\theta_1} n(\theta_1, (\theta_j)_{2 \leq j \leq J})$ . That is, we can increase the diversity of the population  
 126 by increasing the novelty of  $\theta_1$  w.r.t. the list  $(\theta_j)_{2 \leq j \leq J}$ . In practice, we replace this list by a list of  
 127 nearest neighbors of  $\theta_1$ , as this is easier to compute and the elements of  $(\theta_j)_{2 \leq j \leq J}$  tend to be among  
 128 the nearest neighbors of  $\theta_1$ .

129 Under this form, the diversity gradient cannot benefit from the variance reduction methods in the RL  
 130 literature to efficiently compute policy gradients Sutton et al. (1999). To this end, we need to express  
 131 it as a gradient over the expectation of a sum of scalar quantities obtained by policy  $\pi_{\theta_1}$  at each step  
 132 when interacting with the environment. Therefore, to build a DPG, we need information about the  
 133 novelty of a solution at the time step level. To do so, we introduce a novel space  $\mathcal{D}$ , dubbed *state*  
 134 *descriptor space* and a *state descriptor extraction function*  $\psi : \mathcal{S} \rightarrow \mathcal{D}$ . We assume  $\mathcal{D}$  and  $\mathcal{B}$  have the  
 135 same dimension. Similarly to the novelty of a solution, we now define the novelty of a state  $s$  w.r.t.  $J$   
 136 other states  $(s_j)_{j=1,\dots,J}$  as  $n : \mathcal{S} \times \mathcal{S}^J \rightarrow \mathbb{R}$  such that  $n(s, (s_j)_{j=1,\dots,J}) = \sum_{j=1}^J \|\psi(s), \psi(s_j)\|_{\mathcal{D}}$ ,  
 137 where  $\|\cdot\|_{\mathcal{D}}$  is a distance metric over  $\mathcal{D}$ .

138 Now, we need to link novelty defined at the time step level to novelty defined at the solution level. We  
 139 define the novelty of a state w.r.t. a set of solutions. We say that a state is novel w.r.t. some solutions  
 140 if the state is novel w.r.t. to the states visited by these solutions. More formally:

$$n(s, (\theta_j)_{j=1,\dots,J}) = \sum_{j=1}^J \mathbb{E}_{\pi_{\theta_j}} \sum_t \|\psi(s), \psi(s_t)\|_{\mathcal{D}}. \quad (2)$$

141 While we adopt this definition in this paper, one might as well consider other definitions where, for  
 142 instance, a state is compared to states that have been visited at the same time step during another  
 143 episode. In this context, if the following relation is satisfied:

$$\mathbb{E}_{\pi_{\theta_1}} \sum_t n(s_t, (\theta_j)_{2 \leq j \leq J}) = n(\theta_1, (\theta_j)_{2 \leq j \leq J}), \quad (3)$$

144 then we can compute the DPG of  $d$  w.r.t.  $\theta_1$  as

$$\nabla_{\theta_1}^{DPG} = \nabla_{\theta_1} \mathbb{E}_{\pi_{\theta_1}} \sum_t n(s_t, (\theta_j)_{2 \leq j \leq J}). \quad (4)$$

145 This expression corresponds to the classical policy gradient setting where  $\gamma = 1$  and where the  
 146 corresponding reward signal, here dubbed diversity reward, is computed as  $r_t^D = n(s_t, (\theta_j)_{2 \leq j \leq J})$ .  
 147 Therefore, this gradient can be computed using any PG estimation technique replacing the environ-  
 148 ment reward by the diversity reward  $r_t^D$ .

149 Equation (3) enforces a relation between  $\mathcal{B}$  and  $\mathcal{D}$  and between extraction functions  $\psi$  and  $\xi$ . In  
 150 practice, it may be hard to define the behavior descriptor and state descriptor of a solution that satisfy  
 151 this relation while being meaningful to the problem at hand and tractable. But a strict equality is not  
 152 necessary. It suffices that an increase on the left-hand side implies an increase on the right-hand side  
 153 so that we can still update  $\theta_1$  using (4). Furthermore, when this is not the case, the diversity gradient  
 154 update might not result in an increase of diversity in the behavior descriptor space, but in that case the  
 155 MAP-Elites insertion rule will remove the corresponding solution. We show in Section 6 that we can  
 156 define descriptors that do not satisfy the above relation all the time, but still give satisfactory results.

<sup>1</sup>Remark:  $\theta_2$  can appear twice in the list  $(\theta_j)_{2 \leq j \leq J}$

## 157 4 Related Work

158 A distinguishing feature of our approach is that we combine diversity seeking at the level of trajectories  
159 using behavior descriptors and diversity seeking in the state space using state descriptors. The former  
160 is used by MAP-Elites to select solutions from the grid and contributes structural bias towards diversity,  
161 whereas the latter is used during policy gradient steps in the RL part, see Figure 2b. We organize the  
162 literature review below according to this split between two types of diversity seeking mechanisms.

### 163 QD search in the solution space

164 Simultaneously maximizing diversity and performance is the central goal of QD methods (Pugh  
165 et al., 2016; Cully & Demiris, 2017). Among the various possible combinations offered by the  
166 QD framework, Novelty Search with Local Competition (NSLC) (Lehman & Stanley, 2011b) and  
167 MAP-Elites (Mouret & Clune, 2015) are the two most popular algorithms. NSLC builds on the Novelty  
168 Search (NS) algorithm (Lehman & Stanley, 2011a) and maintains an unstructured archive of solutions  
169 selected for their local performance while MAP-Elites uniformly samples individuals from a structured  
170 grid that discretizes the BD space. Not clear in its current form. I suggest: "QD-PG uses the standard  
171 grid of MAP-Elites. However, we also show in Appendix F that QD-PG can be used with alternative  
172 archive structures.

173 With the objective of improving their data-efficiency, QD-ES algorithms that combine QD and ESs,  
174 such as NSR-ES and NSRA-ES, have been applied to challenging continuous control environments in  
175 Conti et al. (2018). But, as outlined in Colas et al. (2020), they suffer from poor sample efficiency  
176 and the diversity and environment reward functions could be mixed in a more efficient way. In that  
177 respect, the most closely related work w.r.t. ours is ME-ES (Colas et al., 2020). The ME-ES algorithm  
178 also optimizes quality and diversity using MAP-Elites and two ES populations. Using these methods  
179 was shown to be critically more efficient than population-based GA algorithms (Salimans et al., 2017),  
180 but our results show that they are still less sample efficient than off-policy deep RL methods, as they  
181 do not leverage the analytical computation of the policy gradient at the time step level. To the best  
182 of our knowledge, no QD or ES algorithm use an explicit critic for both performance and diversity,  
183 resulting in even higher data-efficiency.

### 184 QD search in the state or action spaces

185 Seeking for diversity in the space of states or actions is generally framed into the RL framework. This  
186 is the case of algorithms maintaining a population of RL agents for exploration without an explicit  
187 diversity criterion (Jaderberg et al., 2017) or algorithms explicitly looking for diversity but in the  
188 action space rather than in the state space like ARAC (Doan et al., 2019), P3S-TD3 (Jung et al., 2020)  
189 and DVD (Parker-Holder et al., 2020).

190 An exception is Stanton & Clune (2016) who define a notion of *intra-life novelty* that is similar to  
191 our state novelty defined in Section 3. However, their novelty relies on skills rather than states. Our  
192 work is also related to algorithms using RL mechanisms to search for diversity only (Eysenbach et al.,  
193 2018; Pong et al., 2019; Lee et al., 2019; Islam et al., 2019). These methods have proven useful in  
194 sparse reward situations, but they are inherently limited when the reward signal can orient exploration,  
195 as they ignore it. Other works sequentially combine diversity seeking and RL. The GEP-PG algorithm  
196 Colas et al. (2018) combines a diversity seeking component, namely *Goal Exploration Processes*  
197 (Forestier et al., 2017) and the DDPG deep RL algorithm (Lillicrap et al., 2015). This sequential  
198 combination of exploration-then-exploitation is also present in GO-EXPLORE (Ecoffet et al., 2019).  
199 Again, this approach is limited when the reward signal can help driving the exploration process to  
200 efficient solutions. These sequential approaches first look for diversity in the behavior descriptor  
201 space, then optimize performance in the state action space, whereas we do so simultaneously in the  
202 behavior descriptor space and in the state space.

203 To the best of our knowledge, QD-PG is the first algorithm optimizing both diversity and performance  
204 in the solution and in the state space, using a sample-efficient policy gradient computation method  
205 for the latter.

## 206 5 Methods

207 Our full algorithm is called QD-PG, its pseudo code is given in Appendix A and its architecture is  
208 depicted in Figure 2. QD-PG is an iterative algorithm based on MAP-Elites that replaces random

209 mutations with policy gradient updates. As we consider a continuous action space and want to  
 210 improve sample efficiency by using an off-policy policy gradient method, we rely on the Twin  
 211 Delayed Deterministic Policy Gradient (TD3) algorithm (Fujimoto et al., 2018). See Appendix B for  
 212 a detailed description of TD3.

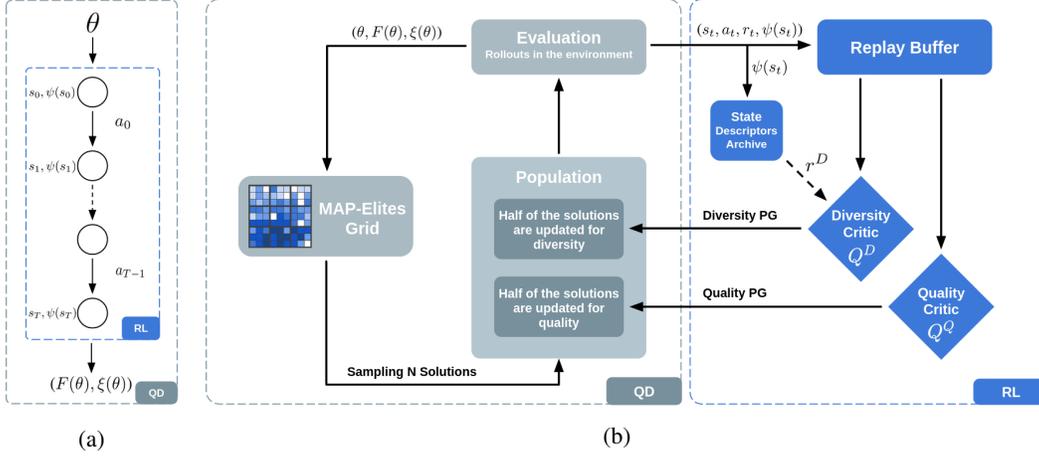


Figure 2: (a): The RL part of QD-PG operates at the time step level while the QD part operates at the controller level, considering the MDP as a black box. (b) One QD-PG iteration consists of three phases: 1) A new population of solutions is sampled from the MAP-Elites grid. 2) These solutions are updated by an off-policy RL agent: half of the solutions are optimized for quality and the other half for diversity. The RL agent leverages one shared critic for each objective. 3) The newly obtained solutions are evaluated in the environment. Transitions are stored in a replay buffer while the updated solutions, their final scores and behavior descriptors are stored in the MAP-Elites grid.

213 QD-PG maintains three permanent structures. In the QD part, a MAP-Elites grid stores the most novel  
 214 and performing solutions. In the RL part, a replay buffer contains all transitions collected when  
 215 evaluating solutions and an archive  $\mathbb{A}$  stores all state descriptors obtained so far. QD-PG starts with an  
 216 initial population of random solutions, evaluates them and inserts them into the MAP-Elites grid. At  
 217 each iteration, solutions are sampled from the grid, copied, and updated. The updated solutions are  
 218 then evaluated through one rollout in the environment and inserted into the grid according to insertion  
 219 rules. Transitions collected during evaluation are stored in the replay buffer, and state descriptors  
 220 are stored in the archive  $\mathbb{A}$ . Note that these state descriptors are first filtered to avoid insertion in the  
 221 archive of multiple state descriptors that are too close to each other.

222 During the update step, half the population is updated with QPG ascent and the other half with DPG  
 223 ascent. The choice of whether an agent is updated for quality or diversity is random, meaning that it  
 224 can be updated for quality and later for diversity if selected again. To justify this design, we show in  
 225 Section 6 that updating consecutively for quality and diversity outperforms updating based on joint  
 226 criteria. Both gradients are computed from batches of transitions sampled from the replay buffer. The  
 227 QPG is computed as usual from rewards whereas for DPG, we get fresh novelty rewards as

$$r_t^D = \sum_{j=1}^J \|\psi(s_t), \psi(s_j)\|_{\mathcal{D}}, \quad (5)$$

228 where  $(s_j)_{j=1, \dots, J}$  are the  $J$  nearest neighbors of state  $s_t$  in the archive  $\mathbb{A}$ . Diversity rewards  
 229 must be recomputed at each update because  $\mathbb{A}$  changes during training. Following Equation (2),  
 230 diversity rewards should be computed as the sum of the distances between the descriptor of  $s_t$  and  
 231 the descriptors of all the states visited by a list of  $J$  solutions. In practice, we consider the  $J$  nearest  
 232 neighbors of  $s_t$ . This choice simplifies the algorithm and is faster and works well in practice.

233 TD3 relies on a parameterized critic to reduce the variance of its policy gradient estimate. In QD-PG,  
 234 we maintain two parameterized critics  $Q_w^D$  and  $Q_v^Q$ , respectively dubbed diversity and quality critics,  
 235 every time a policy gradient is computed, QD-PG also updates the corresponding critic. In fact, as  
 236 in TD3, we use pairs of critics and target critics to fight the overestimation bias. We share the critic  
 237 parameters among the population as in Pourchot & Sigaud (2018). Reasons for doing so come from

238 the fact that diversity is not stationary, as it depends on the current population. If each agent had  
239 its own diversity critic, since an agent may not be selected for a large number of generations before  
240 being selected again, its critic would convey an outdated picture of the evolving diversity. We tried  
241 this solution, and it failed. A side benefit of critic sharing is that both critics become accurate faster as  
242 they combine experience from all agents. Additional details on QD-PG implementation are available  
243 in Appendix C.

## 244 6 Experiments

245 In this section, we intend to answer the following matters: 1. Can QD-PG produce collections of  
246 diverse and high-performing neural policies and what are the advantages to do so? 2. Is QD-PG more  
247 sample efficient than its QD competitors? 3. To what extent are the considered benchmarks difficult  
248 for classical policy gradients methods? 4. What is the usefulness of the different components of  
249 QD-PG?

### 250 Environments

251 We assess QD-PG capabilities in continuous control environments that exhibit high dimensional  
252 observation and action spaces as well as strong exploration difficulties. Two types of reward signals,  
253 dubbed sparse and deceptive, are known to be particularly difficult for classical RL methods. These  
254 rewards appear in many applications such as robotics or combinatorial optimization. Sparse rewards  
255 are obtained if a given condition is specified, leading to a majority of null rewards and to credit  
256 assignment difficulties. Deceptive rewards are dense signals, i.e., they are non-zero at each time step  
257 but can mislead the search process to some local optimum. In such problems, a good approach to the  
258 exploration-exploitation trade-off is essential. The agent should learn when to ignore the reward signal  
259 and explore to avoid local minima and when to follow it to increase its return. Deceptive environments  
260 constitute a natural choice to highlight QD efficiency to balance exploration and exploitation. In this  
261 study, we consider three OpenAI Gym environments based on the MUJoCo physics engine that all  
262 exhibit strong deceptive rewards (illustrated in Appendix 5). Such environments have been widely  
263 used in previous works (Parker-Holder et al., 2020; Colas et al., 2020; Frans et al., 2018; Shi et al.,  
264 2020) for their deceptive nature, a characteristic that is absent of more widespread continuous control  
265 environments like HALFCHEETAH-V2, HOPPER-V2 or still ANT-V2.

266 In the POINT-MAZE environment, an agent represented as a green sphere must find the exit of the  
267 maze depicted in Figure 4a, represented as a red sphere. An observation contains the agent position  
268 at time  $t$ , and an action corresponds to position increments along the  $x$  and  $y$  axes. The reward is  
269 expressed as the negative Euclidean distance between the center of gravity of the agent and the exit  
270 center. The trajectory length cannot exceed 200 steps.

271 The ANT-MAZE environment is modified from OpenAI Gym ANT-V2 (Brockman et al., 2016) and  
272 also used in (Colas et al., 2020; Frans et al., 2018). In ANT-MAZE, a four-legged ant has to reach  
273 a goal zone located in the lower right part of the maze (colored in green in Figure 4b). Its initial  
274 position is sampled in a small circle located in the maze’s extreme bottom left. As in POINT-MAZE,  
275 the reward is expressed as the negative Euclidean distance between the ant and the center of the goal  
276 zone. Maze walls are organized so that following the gradient of the reward function drives the ant  
277 into a dead-end. In ANT-MAZE, the final performance is defined as the maximum reward received  
278 during an episode. The environment is considered solved when an agent obtains a score superior to  
279  $-10$ , corresponding to reaching the goal zone. An episode consists of 3000 time steps, this horizon  
280 is three times larger than in usual MUJoCo environments, making this environment particularly  
281 challenging for RL based methods (Vemula et al., 2019).

282 Finally, the ANT-TRAP environment also derives from ANT-V2 and is inspired from (Colas et al.,  
283 2020; Parker-Holder et al., 2020). In ANT-TRAP, the four-legged ant initially appears in front of a trap  
284 and must bypass it to run as fast as possible in the forward direction (see Figure 4c), as in ANT-V2,  
285 the reward is computed as the ant velocity on the  $x$ -axis. The trap consists of three walls forming a  
286 dead-end directly in front of the ant, leading to a strong deceptive reward. In this environment, the  
287 trajectory length cannot exceed 1000 steps. As opposed to POINT-MAZE and ANT-MAZE, where the  
288 objective is to reach the exit area, there is no unique way to solve ANT-TRAP and we expect a QD  
289 algorithm to generate various effective solutions as depicted in Figure 1.

## 290 Baselines and Ablations

291 QD-PG is compared to three types of methods. First, to answer question 2, we compare QD-PG to a  
292 family of QD baselines, namely ME-ES, NSR-ES, and NSRA-ES (Colas et al., 2020). Appendix E.1  
293 recaps the properties of all these methods. Second, to answer question 3, we compare QD-PG to a  
294 family of policy gradient baselines. Soft Actor Critic (SAC) (Haarnoja et al., 2018) and the Twin  
295 Delayed Deep Deterministic policy gradient (TD3) (Fujimoto et al., 2018) are continuous control  
296 algorithms achieving state-of-the-art results on MUJOCO benchmarks. Random Network Distillation  
297 (RND) (Burda et al., 2018) is a curiosity-driven RL agent (Schulman et al., 2017) which was shown to  
298 perform well in hard exploration settings. CEM-RL (Pourchot & Sigaud, 2018) mixes Cross-Entropy  
299 Methods (CEM) and RL to evolve a population of agents to maximize quality and obtains state-  
300 of-the-art results MUJOCO benchmarks. Finally, to answer question 4, we propose to investigate  
301 the following matters: Can we replace alternating quality and diversity updates by a single update  
302 that optimizes for the sum of both criteria? Are quality gradients updates alone enough to fill the  
303 MAP-Elites grid? Are diversity gradients updates alone enough to do so? Consequently, we consider  
304 the following ablations of QD-PG: QD-PG SUM computes a gradient to optimize the sum of the quality  
305 and diversity rewards, D-PG applies only diversity gradients to the solutions, and Q-PG applies only  
306 quality gradients, but both D-PG and Q-PG still use QD selection (see Appendix E.1).

307 We compare QD-PG to its ablations and RL competitors in all environments and show results in  
308 Table 1a. Detailed results including graphic charts and coverage maps are given in Appendix E and  
309 more details about the evaluation procedure are given in Appendix E.1.

## 310 7 Results

### 311 1. Can QD-PG produce collections of neural policies and what are the advantages to do so?

312 Table 1a presents QD-PG performances. In all environments, our algorithm manages to find working  
313 solutions that avoid local minima and reach the overall objective. In addition to its exploration  
314 capabilities, QD-PG generates collections of high performing solutions in a single run. During the  
315 ANT-TRAP experiment, the final collection of solutions returned by QD-PG contained, among others,  
316 5 solutions that were within a 10% performance margin from the best one. As illustrated in Figure 1,  
317 these agents typically differ in their gaits and preferred trajectories to circumvent the trap.

318 Generating a collection of diverse solutions comes with  
319 the benefit of having a repertoire of diverse solutions that  
320 can be used as alternatives when the MDP changes (Cully  
321 et al., 2015). We show that QD-PG is more robust than  
322 conventional policy gradient methods by changing the re-  
323 ward signal of the ANT-MAZE environment. We replace  
324 the original goal in the bottom right part of the maze (see  
325 Figure 3) with a new randomly located goal in the maze.  
326 Instead of running QD-PG to optimize for this new objec-  
327 tive, we run a Bayesian optimization process to quickly  
328 find a good solution among the ones already stored in the  
329 grid. With a budget of only 20 solutions to be tested during  
330 the Bayesian optimization process, we are able to quickly  
331 recover a good solution for the new objective. We repeat  
332 this experiment 100 times, each time with a different ran-  
333 dom goal, and obtain an average performance of  $-10$  with  
334 a standard deviation of 9. In other words, 20 interaction episodes (corresponding to 60.000 time  
335 steps) suffice for the adaptation process to find a solution that performs well for the new objective  
336 without the need to re-train agents. More detailed results can be found in Appendix E.3.<sup>2</sup>

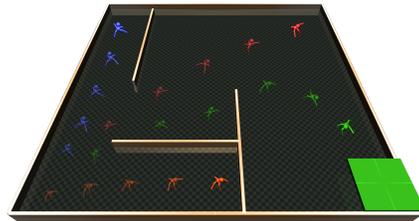


Figure 3: QD-PG produces a collection of diverse solutions. In ANT-MAZE, even after setting new randomly located goals, the MAP-Elites grid still contains solutions that are suited for the new objectives.

### 337 2. Is it more sample efficient than its QD competitors?

338 Table 1b compares QD-PG to Deep Neuroevolution algorithms with a diversity seeking component in  
339 terms of sample efficiency. QD-PG runs on 10 CPU cores for 2 days while its competitors used 1000  
340 CPU cores for the same duration. Nonetheless, QD-PG matches the asymptotic performance of ME-ES  
341 using two orders of magnitude fewer samples, explaining the lower resource requirements.

<sup>2</sup>Videos of QD-PG agents are available at: <https://sites.google.com/view/qd-pg>

Table 1: Results for all environments. **Final Perf.** is the minimum distance to the goal in ANT-MAZE and the episode return in POINT-MAZE and ANT-TRAP. The **Ratio to ours** column compares the sample efficiency of a method to QD-PG.

(a) Comparison to ablations and PG baselines.				(b) Comparison to evolutionary competitors.			
Algorithm	Final Perf. ( $\pm$ std)			Algorithm	ANT-MAZE		
	POINT-MAZE	ANT-MAZE	ANT-TRAP		Final Perf.	Steps to goal	Ratio to ours
QD-PG	-24( $\pm$ 0)	-7( $\pm$ 7)	<b>1541(<math>\pm</math>86)</b>	QD-PG	-7( $\pm$ 7)	<b>1.15e8</b>	<b>1</b>
QD-PG SUM	-25( $\pm$ 1)	-5( $\pm$ 3)	1018( $\pm$ 6)	CEM-RL	-26( $\pm$ 0)	$\infty$	$\infty$
D-PG	-37( $\pm$ 3)	-2( $\pm$ 0)	1016( $\pm$ 8)	ME-ES	-5( $\pm$ 1)	2.4e10	209
Q-PG	-128( $\pm$ 0)	-26( $\pm$ 0)	1175( $\pm$ 79)	NSR-ES	-26( $\pm$ 0)	$\infty$	$\infty$
CEM-RL	-312( $\pm$ 1)	-26( $\pm$ 0)	934( $\pm$ 22)	NSRA-ES	-2( $\pm$ 1)	2.1e10	182
SAC	-127( $\pm$ 1)	-59( $\pm$ 1)	1049( $\pm$ 21)				
TD3	-130( $\pm$ 2)	-26( $\pm$ 0)	1131( $\pm$ 7)				
RND	-35( $\pm$ 10)	-27( $\pm$ 1)	978( $\pm$ 61)				

342 We see three reasons for the improved sample efficiency of QD-PG: 1) QD-PG leverages a replay  
 343 buffer and can re-use each sample several times. 2) QD-PG leverages novelty at the state level and  
 344 can exploit all collected transitions to maximize quality and diversity. For instance, in ANT-MAZE,  
 345 a trajectory brings 3000 samples to QD-PG while standard QD methods would consider it a unique  
 346 sample. 3) PG exploits the analytical gradient between the neural network weights and the resulting  
 347 policy action distribution and estimates only the impact of the distribution on the return. By contrast,  
 348 standard QD methods directly estimate the impact on the return of randomly modifying the weights.

### 349 3. To what extent the considered benchmarks are difficult for policy gradients methods?

350 Table 1a compares QD-PG to state-of-the-art policy gradient algorithms and validates that classical  
 351 policy gradient methods fail to find optimal solutions in deceptive environments. TD3 quickly  
 352 converges to local minima of performance resulting from being attracted in dead-ends by the deceptive  
 353 gradients. While we may expect SAC to better explore due to entropy regularization, it also converges  
 354 to that same local minima in ANT-TRAP and POINT-MAZE. Besides, despite its exploration mechanism  
 355 based on CEM, CEM-RL also quickly converges to local optima in all benchmarks, confirming the  
 356 need for a dedicated diversity seeking component. RND, which adds an exploration bonus used as  
 357 an intrinsic reward (see Appendix G for more details), also demonstrates performances inferior to  
 358 QD-PG in all environments but manages to solve POINT-MAZE. In ANT-MAZE and ANT-TRAP, as  
 359 shown in Appendix G.2, RND extensively explores the BD space but fails to obtain high returns.

### 360 4. What is the usefulness of the different components of QD-PG ?

361 The ablation study in Table 1a shows that when maximising quality only, Q-PG fails due to the  
 362 deceptive nature of the reward and when maximizing diversity only, D-PG sufficiently explores to  
 363 solve the problem in both POINT-MAZE and ANT-MAZE but requires more steps and finds lower-  
 364 performing solutions. When optimizing simultaneously for quality and diversity, QD-PG SUM fails  
 365 to learn in ANT-TRAP and manages to solve the task in ANT-MAZE but requires more samples than  
 366 QD-PG. We hypothesize that quality and diversity rewards may give rise to conflicting gradients. For  
 367 instance, at the beginning of training in ANT-TRAP, the quality reward drives the ant forward whereas  
 368 the diversity reward drives it back to escape the trap and explore the environment. Therefore, both  
 369 rewards cancel each other, preventing any learning. This study validates the usefulness of QD-PG  
 370 components: 1) optimizing for diversity is required to overcome the deceptive nature of the reward;  
 371 2) adding quality optimization provides better asymptotic performance; 3) it is better to disentangle  
 372 quality and diversity updates.

## 373 8 Conclusion

374 This paper is the first to introduce a diversity gradient to explore diversity both at the state and  
 375 skill levels. Based on this component we proposed a novel algorithm, QD-PG, inspired from the  
 376 Quality-Diversity literature, that produces collections of diverse and high-performing neural policies  
 377 in a sample-efficient manner. We showed experimentally that QD-PG generates several solutions  
 378 that achieve high returns in challenging exploration problems. Finally, we demonstrated that in  
 379 a few interactions with the environment, QD-PG finds alternative solutions that still obtain good  
 380 performance when the MDP changes.

## 381 References

- 382 Alvarez, A., Dahlskog, S., Font, J., and Togelius, J. Empowering quality diversity in dungeon design  
383 with interactive constrained map-elites. In *2019 IEEE Conference on Games (CoG)*, pp. 1–8. IEEE,  
384 2019.
- 385 Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W.  
386 Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 387 Burda, Y., Edwards, H., Storkey, A., and Klimov, O. Exploration by random network distillation.  
388 *arXiv preprint arXiv:1810.12894*, 2018.
- 389 Cazenille, L., Bredeche, N., and Aubert-Kato, N. Exploring self-assembling behaviors in a swarm of  
390 bio-micro-robots using surrogate-assisted map-elites. *arXiv preprint arXiv:1910.00230*, 2019.
- 391 Colas, C., Sigaud, O., and Oudeyer, P.-Y. GEP-PG: Decoupling exploration and exploitation in deep  
392 reinforcement learning algorithms. *arXiv preprint arXiv:1802.05054*, 2018.
- 393 Colas, C., Madhavan, V., Huizinga, J., and Clune, J. Scaling map-elites to deep neuroevolution. In  
394 *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*, pp. 67–75, 2020.
- 395 Conti, E., Madhavan, V., Such, F. P., Lehman, J., Stanley, K., and Clune, J. Improving exploration in  
396 evolution strategies for deep reinforcement learning via a population of novelty-seeking agents. In  
397 *Advances in neural information processing systems*, pp. 5027–5038, 2018.
- 398 Cully, A. and Demiris, Y. Quality and diversity optimization: A unifying modular framework. *IEEE*  
399 *Transactions on Evolutionary Computation*, 22(2):245–259, 2017.
- 400 Cully, A., Clune, J., Tarapore, D., and Mouret, J.-B. Robots that can adapt like animals. *Nature*, 521  
401 (7553):503–507, 2015.
- 402 Doan, T., Mazouze, B., Durand, A., Pineau, J., and Hjelm, R. D. Attraction-repulsion actor-critic for  
403 continuous control reinforcement learning. *arXiv preprint arXiv:1909.07543*, 2019.
- 404 Doncieux, S., Laflaquière, A., and Coninx, A. Novelty search: a theoretical perspective. In  
405 *Proceedings of the Genetic and Evolutionary Computation Conference*, pp. 99–106, 2019.
- 406 Ecoffet, A., Huizinga, J., Lehman, J., Stanley, K. O., and Clune, J. Go-explore: a new approach for  
407 hard-exploration problems. *arXiv preprint arXiv:1901.10995*, 2019.
- 408 Eysenbach, B., Gupta, A., Ibarz, J., and Levine, S. Diversity is all you need: Learning skills without  
409 a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- 410 Forestier, S., Mollard, Y., and Oudeyer, P.-Y. Intrinsically motivated goal exploration processes with  
411 automatic curriculum learning. *arXiv preprint arXiv:1708.02190*, 2017.
- 412 Frans, K., Ho, J., Chen, X., Abbeel, P., and Schulman, J. Meta learning shared hierarchies. *Proc. of*  
413 *ICLR*, 2018.
- 414 Fujimoto, S., Van Hoof, H., and Meger, D. Addressing function approximation error in actor-critic  
415 methods. *arXiv preprint arXiv:1802.09477*, 2018.
- 416 Haarnoja, T., Zhou, A., Hartikainen, K., Tucker, G., Ha, S., Tan, J., Kumar, V., Zhu, H., Gupta, A.,  
417 Abbeel, P., et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*,  
418 2018.
- 419 Islam, R., Ahmed, Z., and Precup, D. Marginalized state distribution entropy regularization in policy  
420 optimization. *arXiv preprint arXiv:1912.05128*, 2019.
- 421 Jaderberg, M., Dalibard, V., Osindero, S., Czarnecki, W. M., Donahue, J., Razavi, A., Vinyals, O.,  
422 Green, T., Dunning, I., Simonyan, K., et al. Population-based training of neural networks. *arXiv*  
423 *preprint arXiv:1711.09846*, 2017.
- 424 Jung, W., Park, G., and Sung, Y. Population-guided parallel policy search for reinforcement learning.  
425 In *International Conference on Learning Representations*, 2020.

- 426 Khadka, S. and Tumer, K. Evolution-guided policy gradient in reinforcement learning. In *Neural*  
427 *Information Processing Systems*, 2018.
- 428 Khadka, S., Majumdar, S., Miret, S., Tumer, E., Nassar, T., Dwiel, Z., Liu, Y., and Tumer, K.  
429 Collaborative evolutionary reinforcement learning. *arXiv preprint arXiv:1905.00976*, 2019.
- 430 Koos, S., Mouret, J.-B., and Doncieux, S. The transferability approach: Crossing the reality gap in  
431 evolutionary robotics. *IEEE Transactions on Evolutionary Computation*, 17(1):122–145, 2012.
- 432 Lee, L., Eysenbach, B., Parisotto, E., Xing, E., Levine, S., and Salakhutdinov, R. Efficient exploration  
433 via state marginal matching. *arXiv preprint arXiv:1906.05274*, 2019.
- 434 Lehman, J. and Stanley, K. O. Abandoning objectives: Evolution through the search for novelty  
435 alone. *Evolutionary computation*, 19(2):189–223, 2011a.
- 436 Lehman, J. and Stanley, K. O. Evolving a diversity of virtual creatures through novelty search and  
437 local competition. In *Proceedings of the 13th annual conference on Genetic and evolutionary*  
438 *computation*, pp. 211–218, 2011b.
- 439 Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D.  
440 Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- 441 Mouret, J.-B. and Clune, J. Illuminating search spaces by mapping elites. *arXiv preprint*  
442 *arXiv:1504.04909*, 2015.
- 443 Nasiriany, S., Pong, V. H., Lin, S., and Levine, S. Planning with goal-conditioned policies. *arXiv*  
444 *preprint arXiv:1911.08453*, 2019.
- 445 Parker-Holder, J., Pacchiano, A., Choromanski, K., and Roberts, S. Effective diversity in population-  
446 based reinforcement learning. In *Neural Information Processing Systems*, 2020.
- 447 Pong, V. H., Dalal, M., Lin, S., Nair, A., Bahl, S., and Levine, S. Skew-fit: State-covering self-  
448 supervised reinforcement learning. *arXiv preprint arXiv:1903.03698*, 2019.
- 449 Pourchot, A. and Sigaud, O. Cem-rl: Combining evolutionary and gradient-based methods for policy  
450 search. *arXiv preprint arXiv:1810.01222*, 2018.
- 451 Pugh, J. K., Soros, L. B., and Stanley, K. O. Quality diversity: A new frontier for evolutionary  
452 computation. *Frontiers in Robotics and AI*, 3:40, 2016.
- 453 Salimans, T., Ho, J., Chen, X., Sidor, S., and Sutskever, I. Evolution strategies as a scalable alternative  
454 to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- 455 Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. High-dimensional continuous control  
456 using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- 457 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization  
458 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 459 Shi, L., Li, S., Zheng, Q., Yao, M., and Pan, G. Efficient novelty search through deep reinforcement  
460 learning. *IEEE Access*, 8:128809–128818, 2020.
- 461 Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., and Riedmiller, M. Deterministic policy  
462 gradient algorithms. In *Proceedings of the 30th International Conference in Machine Learning*,  
463 2014.
- 464 Stanton, C. and Clune, J. Curiosity search: producing generalists by encouraging individuals to  
465 continually explore and acquire skills throughout their lifetime. *PloS one*, 11(9):e0162235, 2016.
- 466 Sutton, R. S., McAllester, D. A., Singh, S. P., Mansour, Y., et al. Policy gradient methods for  
467 reinforcement learning with function approximation. In *NIPs*, volume 99, pp. 1057–1063. Citeseer,  
468 1999.
- 469 Vemula, A., Sun, W., and Bagnell, J. Contrasting exploration in parameter and action space: A zeroth-  
470 order optimization perspective. In *The 22nd International Conference on Artificial Intelligence*  
471 *and Statistics*, pp. 2926–2935. PMLR, 2019.

472 **Checklist**

- 473 1. For all authors...
- 474 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
475 contributions and scope? [Yes]
- 476 (b) Did you describe the limitations of your work? [Yes]
- 477 (c) Did you discuss any potential negative societal impacts of your work? [No] **We believe**  
478 **that this work, in itself, is not prone to have any negative societal impact.**
- 479 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
480 them? [Yes]
- 481 2. If you are including theoretical results...
- 482 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 483 (b) Did you include complete proofs of all theoretical results? [N/A]
- 484 3. If you ran experiments...
- 485 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
486 mental results (either in the supplemental material or as a URL)? [Yes] **The code and**  
487 **instructions to run it are available in the supplementary materials.**
- 488 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were  
489 chosen)? [Yes] **Implementation details, hardware details and hyperparameters**  
490 **are presented in Appendix C.**
- 491 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
492 ments multiple times)? [Yes] **Yes, we report mean and variance for all experiments,**  
493 **both graphically and in result tables.**
- 494 (d) Did you include the total amount of compute and the type of resources used (e.g.,  
495 type of GPUs, internal cluster, or cloud provider)? [Yes] **Computational details are**  
496 **provided in Appedix C.**
- 497 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 498 (a) If your work uses existing assets, did you cite the creators? [Yes] **We use open**  
499 **sourced RL environments.**
- 500 (b) Did you mention the license of the assets? [Yes] **We cite the Mujoco physics engine,**  
501 **for which we have licenses.**
- 502 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
503 **We provide appendices, source code and a demonstration website.**
- 504 (d) Did you discuss whether and how consent was obtained from people whose data you're  
505 using/curating? [N/A]
- 506 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
507 information or offensive content? [N/A]
- 508 5. If you used crowdsourcing or conducted research with human subjects...
- 509 (a) Did you include the full text of instructions given to participants and screenshots, if  
510 applicable? [N/A] **We did not used crowdsourcing or conducted research with**  
511 **human subjects.**
- 512 (b) Did you describe any potential participant risks, with links to Institutional Review  
513 Board (IRB) approvals, if applicable? [N/A] **This work did not involve research**  
514 **with human subjects**
- 515 (c) Did you include the estimated hourly wage paid to participants and the total amount  
516 spent on participant compensation? [N/A] **This work did not involve research with**  
517 **human subjects**