

Synergizing Quality-Diversity with Descriptor-Conditioned Reinforcement Learning

ANONYMOUS AUTHOR(S)

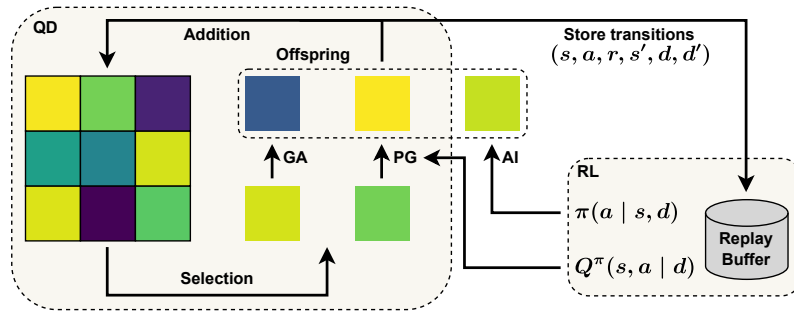


Fig. 1. DCG-MAP-ELITES-AI implements a conventional MAP-ELITES loop comprising selection, variation, evaluation, addition and leverages two complementary variation operators: a standard Genetic Algorithm (GA) variation operator for diversity and a descriptor-conditioned Policy Gradient (PG) variation operator for quality. Concurrently to the critic’s training, the knowledge of the archive is distilled in the descriptor-conditioned actor. In turn, this versatile actor is injected (AI) in the offsprings at each iteration.

A fundamental trait of intelligence involves finding novel and creative solutions to address a given challenge or to adapt to unforeseen situations. Reflecting this, Quality-Diversity optimization is a family of Evolutionary Algorithms, that generates collections of both diverse and high-performing solutions. Among these, MAP-ELITES is a prominent example, that has been successfully applied to a variety of domains, including evolutionary robotics. However, MAP-ELITES performs a divergent search with random mutations originating from Genetic Algorithms, and thus, is limited to evolving populations of low-dimensional solutions. PGA-MAP-ELITES overcomes this limitation using a gradient-based variation operator inspired by deep reinforcement learning which enables the evolution of large neural networks. Although high-performing in many environments, PGA-MAP-ELITES fails on several tasks where the convergent search of the gradient-based variation operator hinders diversity. In this work, we present three contributions: (1) we enhance the Policy Gradient variation operator with a descriptor-conditioned critic that reconciles diversity search with gradient-based methods, (2) we leverage the actor-critic training to learn a descriptor-conditioned policy at no additional cost, distilling the knowledge of the population into one single versatile policy that can execute a diversity of behaviors, (3) we exploit the descriptor-conditioned actor by injecting it in the population, despite network architecture differences. Our method, DCG-MAP-ELITES-AI, achieves equal or higher QD score and coverage compared to all baselines on seven challenging continuous control locomotion tasks.

CCS Concepts: • **Computing methodologies** → **Evolutionary robotics; Sequential decision making.**

Additional Key Words and Phrases: Quality-Diversity, Reinforcement Learning, Neuroevolution, MAP-Elites, Policy Gradient

ACM Reference Format:

Anonymous Author(s). 2018. Synergizing Quality-Diversity with Descriptor-Conditioned Reinforcement Learning. In . ACM, New York, NY, USA, 30 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

A fascinating aspect of evolution is its ability to generate a variety of different species, each being adapted to their niche. Inspired by this idea, Quality-Diversity (QD) optimization is a family of evolutionary algorithms that aims to generate a set of both high-performing and diverse solutions to a single problem [5, 9, 35]. Contrary to traditional optimization methods that return a single high-performing solution, the goal of QD algorithms is to illuminate a search space of interest called *descriptor space* [30]. Producing a large collection of diverse and effective solutions enables to get multiple alternatives to solve a single problem, which is useful in robotics to improve robustness, recover from damage [8] or reduce the reality gap [6]. Furthermore, conventional optimization methods are prone to get stuck in local optima, whereas keeping a repertoire of diverse solutions to a given problem can help to find stepping stones that lead to globally better solutions [30, 31]. Another benefit of diversity search is efficient exploration in problems where the reward signal is sparse or deceptive [4, 10, 34].

MAP-ELITES [30] is a conceptually simple but effective QD optimization algorithm that has shown competitive results in a variety of applications, to generate large collections of diverse skills. However, MAP-ELITES relies on random variations that can cause slow convergence in large search spaces [7, 31, 34], making it inadequate to evolve neural networks with a large number of parameters.

In contrast, Deep Reinforcement Learning (RL) [29] algorithms combine reinforcement learning with the directed search power of gradient-based methods in order to learn a single optimal solution. RL has led to remarkable accomplishments in various areas, including in discrete environments like video games [45], board games [39] and in continuous control domains for locomotion [21, 23] and manipulation [32]. These achievements highlight the exceptional capabilities of RL algorithms in addressing specific challenges. Especially, policy gradient methods have shown state-of-the-art results in learning large neural network policies with thousands of parameters in high-dimensional and continuous domains [21, 28, 40].

PGA-MAP-ELITES [31] is an extension of MAP-ELITES that integrates the sample efficiency of RL algorithms using TD3 [19]. It combines a Policy Gradient (PG) variation operator for efficient fitness improvement, coupled with the usual Genetic Algorithm (GA) variation operator. The PG variation operator leverages gradients derived from RL to drive mutations towards the global fitness optimum and is supported by the divergent search of the GA variation operator for both exploration and optimization [13]. Other recent works have also introduced methods to combine the strength of QD algorithms with reinforcement learning [34, 42] on complex robotics tasks.

PGA-MAP-ELITES achieves state-of-the-art performances in most of the environments considered so far in the literature [31, 34, 42]. However, the PG variation operator becomes ineffective in tasks where the global optimum is in an area of the search space that is not likely to produce offspring that are added to the archive. For example, consider a locomotion task where the fitness is the opposite of the energy consumption and the descriptor is defined as the final position of the robot. The global optimum for the fitness is the solution that does not move in order to minimize energy consumption. Thus, the PG variation operator will encourage solutions to stay motionless, collapsing their descriptors to a single point, the descriptor of the global optimum. Consequently, the PG variation operator generates offspring that are discarded and no interesting stepping stone is found, thereby hindering diversity.

DCG-MAP-ELITES GECCO [12] builds upon PGA-MAP-ELITES algorithm by enhancing the PG variation operator with a descriptor-conditioned critic that provides gradients depending on a target descriptor. The descriptor-conditioned critic takes as input a state and a target descriptor to evaluate actions. Thus, the PG variation operator can mutate

solutions to produce offsprings with higher fitness while targeting a desired descriptor, thereby avoiding to collapse their descriptors to a single point.

Furthermore, the descriptor-conditioned critic undergoes training utilizing the RL algorithm TD3 that requires to train an actor in parallel. We take advantage of this intertwined actor-critic training to make the actor ‘descriptor-conditioned’ as well, allowing it to take actions based not only on the current state but also on a target descriptor we want to achieve. Thus, instead of taking actions that maximize the fitness globally, the actor now takes actions that maximize the fitness while achieving a target descriptor. At the end of training, the result is a versatile agent that can achieve the diversity of behaviors contained in the archive while obtaining similar fitness performance, negating the burden of dealing with a collection of thousands of solutions. In addition to archive distillation, DCG-MAP-ELITES GECCO has been shown to improve performance significantly over PGA-MAP-ELITES on omnidirectional tasks, while maintaining similar performance on unidirectional tasks where no improvement was expected.

Finally, drawing inspiration from PGA-MAP-ELITES that injects the actor in the population at each generation, we extend the original DCG-MAP-ELITES GECCO version [12] with a descriptor-conditioned Actor Injection (AI), that enables to inject the versatile actor in the population, despite network architecture differences.

In summary, we introduce DCG-MAP-ELITES-AI (Descriptor-Conditioned Gradients MAP-Elites with Actor Injection) that extends DCG-MAP-ELITES GECCO and present three contributions: (1) we enhance the PG variation operator with a descriptor-conditioned critic, (2) we distill the knowledge of the archive into one single versatile policy at no additional cost, (3) we take advantage of this high-performing and versatile policy to improve the population during training with actor injection, further improving our method. We compare our algorithm to four state-of-the-art QD algorithms on seven challenging continuous control locomotion tasks. Our method, DCG-MAP-ELITES-AI, achieves equal or higher QD score and coverage compared to all baselines on seven challenging continuous control locomotion tasks.

2 BACKGROUND

2.1 Problem Statement

We consider an agent sequentially interacting with an environment at discrete time steps t for an episode of length T . At each time step t , the agent observes a state s_t , takes an action a_t and receives a scalar reward r_t . We model it as a Markov Decision Process (MDP) which comprises a *state space* \mathcal{S} , a continuous *action space* \mathcal{A} , a stationary *transition dynamics distribution* $p(s_{t+1} | s_t, a_t)$ and a *reward function* $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. In this work, a *policy* (also called *solution*) is a deterministic neural network parameterized by $\phi \in \Phi$, and denoted $\pi_\phi: \mathcal{S} \rightarrow \mathcal{A}$. The agent uses its policy to select actions and interact with the environment to give a trajectory of states, actions and rewards. The *fitness* of a solution is given by $F: \Phi \rightarrow \mathbb{R}$, defined as the expected discounted return $\mathbb{E}_{\pi_\phi} [\sum_{t=0}^{T-1} \gamma^t r_t]$.

In this setting, the objective of QD algorithms is to find the highest fitness solutions in each point of the *descriptor space* \mathcal{D} . The descriptor function $D: \Phi \rightarrow \mathcal{D}$ is generally defined by the user and characterizes solutions in a meaningful way for the type of diversity desired. With this notation, our objective is to evolve a population of solutions that are both high-performing with respect to F and diverse with respect to D .

2.2 MAP-ELITES

Multi-dimensional Archive of Phenotypic Elites (MAP-ELITES) [30] is a simple yet effective QD algorithm, that discretizes the descriptor space \mathcal{D} into a multi-dimensional grid of cells called archive \mathcal{X} and searches for the best solution in each

cell, see Algorithm 14. The goal of the algorithm is to return an archive that is filled as much as possible with high-fitness solutions. MAP-ELITES starts by generating random solutions and adding them to the archive. The algorithm then repeats the following steps until a budget of I solutions have been evaluated: (1) a batch of solutions from the archive are uniformly selected and modified through mutations and/or crossovers to produce offspring, (2) the fitnesses and descriptors of the offspring are evaluated, and each offspring is placed in its corresponding cell if and only if the cell is empty or if the offspring has a better fitness than the current solution in that cell, in which case the current solution is replaced. As most evolutionary methods, MAP-ELITES relies on undirected updates that are agnostic to the fitness objective. With a Genetic Algorithm (GA) variation operator, MAP-ELITES performs a divergent search that may cause slow convergence in high-dimensional problems due to a lack of directed search power, and thus, is performing best on low-dimensional search space [31].

2.3 Deep Reinforcement Learning

Deep Reinforcement Learning (RL) [29] combines the reinforcement learning framework with the function approximation capabilities of deep neural networks to represent policies and value functions in high-dimensional state and action spaces. In opposition to black-box optimization methods like evolutionary algorithms, RL leverages the structure of the MDP in the form of the Bellman equation to achieve better sample efficiency. The objective is to find an optimal policy π_ϕ , which maximizes the expected return or fitness $F(\pi_\phi)$. In reinforcement learning, many approaches try to estimate the action-value function $Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{i=0}^{T-t-1} \gamma^i r_{t+i} \mid s_t = s, a_t = a \right]$ defined as the expected discounted return starting from state s , taking action a and thereafter following policy π .

The Twin Delayed Deep Deterministic Policy Gradient (TD3) algorithm [19] is an actor-critic, off-policy reinforcement learning method that achieves state-of-the-art results in environments with large and continuous action space. TD3 indirectly learns a policy π_ϕ via maximization of the action-value function $Q_\theta(s, a)$. The approach is closely connected to Q -learning [19] and tries to approximate the optimal action-value function $Q^*(s, a)$ in order to find the optimal action $\pi^*(s) = \arg \max_a Q^*(s, a)$. However, computing the maximum over action in $\max_a Q_\theta(s, a)$ is intractable in continuous action space, hence it is approximated using $\max_a Q_\theta(s, a) = Q_\theta(s, \pi_\phi(s))$. In TD3, the policy π_ϕ takes actions in the environment and the transitions are stored in a replay buffer. The collected experience is then used to train a pair of critics $Q_{\theta_1}, Q_{\theta_2}$ using temporal difference. Target networks $Q_{\theta_1}', Q_{\theta_2}'$ are updated to slowly track the main networks. Both critics use a single regression target y , calculated using whichever of the two target critics gives a smaller estimated value and using target policy smoothing by sampling a noise $\epsilon \sim \text{clip}(\mathcal{N}(0, \sigma), -c, c)$:

$$y = r(s_t, a_t) + \gamma \min_{i=1,2} Q_{\theta_i}'(s_{t+1}, \pi_{\phi'}(s_{t+1})) + \epsilon \quad (1)$$

Both critics are learned by regression to this target and the policy is learned with a delay, only updated every Δ iterations simply by maximizing Q_{θ_1} with $\max_\phi \mathbb{E} [Q_{\theta_1}(s, \pi_\phi(s))]$. The actor is updated using the deterministic policy gradient:

$$\nabla_\phi J(\phi) = \mathbb{E} \left[\nabla_\phi \pi_\phi(s) \nabla_a Q_{\theta_1}(s, a) \Big|_{a=\pi_\phi(s)} \right] \quad (2)$$

2.4 PGA-MAP-ELITES

Policy Gradient Assisted MAP-Elites (PGA-MAP-ELITES) [31] is an extension of MAP-ELITES that is designed to evolve deep neural networks by combining the directed search power and sample efficiency of RL methods with the exploration capabilities of genetic algorithms, see Algorithm 9. The algorithm follows the usual MAP-ELITES loop of selection, variation, evaluation and addition for a budget of I iterations, but uses two parallel variation operators: half of the

209 offspring are generated using a standard Genetic Algorithm (GA) variation operator and half of the offspring are
210 generated using a Policy Gradient (PG) variation operator. During each iteration of the loop, PGA-MAP-ELITES stores
211 the transitions from offspring evaluation in a replay buffer \mathcal{B} and uses it to train a pair of critics based on the TD3
212 algorithm, described in Algorithm 10. The trained critic is then used in the PG variation operator to update the selected
213 solutions from the archive for m gradient steps to select actions that maximize the approximated action-value function,
214 as described in Algorithm 11. At each iteration, the critics are trained for n steps of gradients descents towards the
215 target described in Equation (1), averaged over N transitions of experience sampled uniformly from the replay buffer \mathcal{B} .
216 The actor learns with a delay Δ via maximization of the critic according to Equation (2).
217
218
219

220 3 RELATED WORK

221 3.1 Scaling QD to Neuroevolution

224 The challenge of evolving diverse solutions in a high-dimensional search space has been an active research subject
225 over recent years. MAP-ELITES-ES [7] scales MAP-ELITES to high-dimensional solutions parameterized by large neural
226 networks. This algorithm leverages Evolution Strategies [36] (ES) to perform a directed search that is more efficient
227 than random mutations used in Genetic Algorithms. Fitness and novelty gradients are estimated locally from many
228 perturbed versions of the parent solution to generate a new one. The population tends towards regions of the parameter
229 space with higher fitness or novelty but it requires to sample and evaluate a large number of solutions, making it
230 particularly data inefficient. To improve sample efficiency, methods that combine MAP-ELITES with RL [31, 33, 34, 42]
231 have emerged and use time step level information to efficiently evolve populations of high-performing and diverse
232 neural network for complex tasks. PGA-MAP-ELITES [31] uses policy gradients for part of its mutations, see Section 2.4
233 for details. CMA-MEGA [42] estimates descriptor gradients with ES and combines the fitness gradient and the descriptor
234 gradients with a CMA-ES mechanism [16, 22]. QD-PG [34] introduces a diversity reward based on the novelty of the
235 states visited and derives a policy gradient for the maximization of those diversity rewards which helps exploration in
236 settings where the reward is sparse or deceptive. PBT-MAP-ELITES [33] mixes MAP-ELITES with a population based
237 training process [25] to optimize hyper-parameters of diverse RL agents. Interestingly, recent work [41] scales the
238 algorithm CMA-MAE [17] to high-dimensional policies on robotics tasks with pure ES while showing comparable data
239 efficiency to QD-RL approaches, but is still outperformed by PGA-MAP-ELITES.
240
241
242
243
244
245
246

247 3.2 Conditioning the critic

248 None of the methods described in the previous section take a descriptor into account when deriving policy gradients used
249 to mutate solutions. In other words, they do not use descriptor-conditioned policies nor descriptor-conditioned critics as
250 our method does. The concept of descriptor-conditioned critic is related to Universal Value Function Approximators [37],
251 extensively used in skill discovery reinforcement learning, a field that share a similar motivation to QD [2]. In VIC,
252 DIAYN, DADS, SMERL [11, 20, 27, 38], the actors and critics are conditioned on a sampled prior but does not correspond
253 to a real posterior like in DCG-MAP-ELITES-AI. Furthermore, those methods use a notion of diversity defined at the
254 step-level rather than trajectory-level like DCG-MAP-ELITES-AI. Moreover, they do not use an archive to store a
255 population, resulting in much smaller sets of final policies. Finally, it has been shown that QD methods are competitive
256 with skill discovery reinforcement learning algorithms [2], specifically for adaptation and hierarchical learning.
257
258
259
260

261 3.3 Archive distillation

262 Distilling the knowledge of an archive into a single policy is an alluring process that reduces the number of parameters
 263 outputted by the algorithm and enables generalization and interpolation/extrapolation. Although distillation is usually
 264 referring to policy distillation – learning the observation/action mapping from a teacher policy – we present archive
 265 distillation as a general term referring to any kind of knowledge transfer from an archive to another model, should it be
 266 the policies, transitions experienced in the environment, full trajectories or discovered descriptors.
 267

268 To the best of our knowledge, only two QD-related works use the concept of archive distillation. Go-Explore [10]
 269 keeps an archive of states and trains a goal-conditioned policy to reproduce the trajectory of the policy that reached
 270 that state. Another related approach is to learn a generative policy network [26] over the policies contained in the
 271 archive. Our approach DCG-MAP-ELITES-AI distills the experience of the archive into a single versatile policy.
 272
 273

274 4 METHODS

275 Algorithm 1 DCG-MAP-ELITES-AI

276 **Require:** GA batch size b_{GA} , PG batch size b_{PG} , Actor Injection batch size b_{AI} , total batch size $b = b_{GA} + b_{PG} + b_{AI}$

277 Initialize archive \mathcal{X} with b random solutions and replay buffer \mathcal{B}

278 Initialize critic networks $Q_{\theta_1}, Q_{\theta_2}$ and actor network π_{ϕ}

279 $i \leftarrow 0$

280 **while** $i < I$ **do**

281 TRAIN_ACTOR_CRITIC($\pi_{\phi}, Q_{\theta_1}, Q_{\theta_2}, \mathcal{B}$)

282 $\pi_{\psi_1}, \dots, \pi_{\psi_b} \leftarrow \text{SELECTION}(\mathcal{X})$

283 $\pi_{\hat{\psi}_1}, \dots, \pi_{\hat{\psi}_{b_{GA}}} \leftarrow \text{VARIATION_GA}(\pi_{\psi_1}, \dots, \pi_{\psi_{b_{GA}}})$

284 $\pi_{\hat{\psi}_{b_{GA}+1}}, \dots, \pi_{\hat{\psi}_{b_{GA}+b_{PG}}} \leftarrow \text{VARIATION_PG}(\pi_{\psi_{b_{GA}+1}}, \dots, \pi_{\psi_{b_{GA}+b_{PG}}}, Q_{\theta_1}, \mathcal{B})$

285 $\pi_{\hat{\psi}_{b_{GA}+b_{PG}+1}}, \dots, \pi_{\hat{\psi}_b} \leftarrow \text{ACTOR_INJECTION}(\pi_{\phi})$

286 ADDITION($\pi_{\hat{\psi}_1}, \dots, \pi_{\hat{\psi}_b}, \mathcal{X}, \mathcal{B}$)

287 $i \leftarrow i + b$

288 **function** ADDITION($\pi_{\hat{\psi}} \dots, \mathcal{X}, \mathcal{B}$)

289 **for** $\pi_{\hat{\psi}} \dots$ **do**

290 $(f, \text{transitions}) \leftarrow F(\pi_{\hat{\psi}}), d \leftarrow D(\pi_{\hat{\psi}})$

291 INSERT($\mathcal{B}, \text{transitions}$)

292 **if** $\mathcal{X}(d) = \emptyset$ or $F(\mathcal{X}(d)) < f$ **then**

293 $\mathcal{X}(d) \leftarrow \pi_{\hat{\psi}}$

294 Our method Descriptor-Conditioned Gradients MAP-Elites with Actor Injection (DCG-MAP-ELITES-AI) overcomes
 295 the limitations of PGA-MAP-ELITES by leveraging a descriptor-conditioned critic to improve the PG variation operator
 296 and concurrently distills the knowledge of the archive in a single versatile policy as a by-product of the actor-critic
 297 training. The pseudocode is provided in Algorithm 1. The algorithm follows the usual MAP-ELITES loop of selection,
 298 variation, evaluation and addition for a budget of I iterations. Two complementary and independent variation operators
 299 are used in parallel: (1) a standard GA operator (2) a descriptor-conditioned PG operator. At each iteration, the transitions
 300 from the evaluation step are stored in a replay buffer and used to train an actor-critic pair based on TD3.
 301
 302

303 Contrary to PGA-MAP-ELITES, the actor-critic pair is descriptor-conditioned. In addition to the state s and action a ,
 304 the critic Q_{θ} ($s, a \mid d$) also depends on the descriptor d and estimates the expected discounted return starting from state
 305
 306
 307
 308
 309
 310
 311
 312

313 s , taking action a and thereafter following policy π and achieving descriptor d . In this work, to achieve descriptor d
 314 means that the trajectory generated by the policy π has descriptor d . In addition to the state s , the actor $\pi_\phi(s | d)$ also
 315 depends on a target descriptor d and maximizes the expected discounted return conditioned on achieving the target
 316 descriptor d . Thus, the goal of the descriptor-conditioned actor is to achieve the desired descriptor d while maximizing
 317 fitness.
 318
 319

321 4.1 Descriptor-Conditioned Critic

322
 323 Instead of estimating the action-value function with $Q_\theta(s, a)$, we want to estimate the descriptor-conditioned action-
 324 value function with $Q_\theta(s, a | d)$. When a policy π interacts with the environment, it generates a trajectory, which is a
 325 sequence of transitions (s, a, r, s') with descriptor d . We extend the definition of a transition (s, a, r, s') to include the
 326 observed descriptor d of the trajectory (s, a, r, s', d) . However, the descriptor is only available at the end of the episode,
 327 therefore the transitions can only be augmented with the descriptor after the episode is completed. In all the tasks we
 328 consider, the reward function is positive $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^+$ and hence, the fitness function F and action-value function
 329 are positive as well. Thus, for any target descriptor $d' \in \mathcal{D}$, we define the descriptor-conditioned critic as equal to the
 330 normal action-value function when the policy achieves the target descriptor d' and as equal to zero when the policy
 331 does not achieve the target descriptor d' . Given a transition (s, a, r, s', d) , and a target descriptor d' sampled in \mathcal{D} ,
 332
 333

$$334 \quad Q_\theta(s, a | d') := \begin{cases} Q_\theta(s, a), & \text{if } d = d' \\ 0, & \text{if } d \neq d' \end{cases} \quad (3)$$

338 However, with this piecewise definition, the descriptor-conditioned action-value function is not continuous and violates
 339 the universal approximation theorem continuity hypothesis [24]. To address this issue, we introduce a similarity function
 340 $S: \mathcal{D}^2 \rightarrow]0, 1]$ defined as $S(d, d') = e^{-\frac{\|d-d'\|_{\mathcal{D}}}{t}}$ to smooth the descriptor-conditioned critic and relax Equation (3) into:
 341

$$342 \quad \begin{aligned} 343 \quad Q_\theta(s, a | d') &= S(d, d') Q_\theta(s, a) = S(d, d') \mathbb{E}_\pi \left[\sum_{i=0}^{T-t-1} \gamma^i r_{t+i} \middle| s, a \right] \\ 344 &= \mathbb{E}_\pi \left[\sum_{i=0}^{T-t-1} \gamma^i S(d, d') r_{t+i} \middle| s, a \right] \end{aligned} \quad (4)$$

348 With Equation (4), we demonstrate that learning the descriptor-conditioned critic is equivalent to scaling the reward by
 349 the similarity $S(d, d')$ between the descriptor of the trajectory d and the target descriptor d' . Therefore, the critic target
 350 in Equation (1) is modified to include the similarity scaling and the descriptor-conditioned actor:
 351

$$352 \quad y = S(d, d') r(s_t, a_t) + \gamma \min_{i=1,2} Q_{\theta_i'}(s_{t+1}, \pi_{\phi_i'}(s_{t+1} | d')) + \epsilon | d' \quad (5)$$

355 If the target descriptor d' is approximately equal to the observed descriptor d of the trajectory $d \approx d'$, then we have
 356 $S(d, d') \approx 1$ so the reward is unchanged. However, if the descriptor d' is different from the observed descriptor d , then
 357 the reward is scaled down to $S(d, d') r(s_t, a_t) \approx 0$. The scaling ensures that the magnitude of the reward depends not
 358 only on the quality of the action a with regards to the fitness function F , but also on achieving the target descriptor
 359 d' . Given one transition (s, a, r, s', d) , we can generate infinitely many critic updates by sampling a target descriptor
 360 $d' \in \mathcal{D}$. This is leveraged in the new actor-critic training introduced with DCG-MAP-ELITES-AI, which is detailed in
 361 Algorithm 2 and Section 4.3.
 362
 363
 364

4.2 Descriptor-Conditioned Actor and Archive Distillation

The training of the critic requires to train an actor π_ϕ to approximate the optimal action a^* , as explained in Section 2.3. However, in this work, the action-value function estimated by the critic is conditioned on a descriptor d . Hence, we don't want π_ϕ to estimate the best action globally, but rather the best action given that it achieves the target descriptor d . Therefore, the actor is extended to a descriptor-conditioned policy $\pi_\phi(s | d)$, that maximizes the descriptor-conditioned critic's value with $\max_\phi \mathbb{E} [Q_\theta(s, \pi_\phi(s | d) | d)]$. The actor is updated using the deterministic policy gradient, see Algorithm 2:

$$\nabla_\phi J(\phi) = \frac{1}{N} \sum \nabla_\phi \pi_\phi(s | d') \nabla_a Q_{\theta_1}(s, a | d')|_{a=\pi_\phi(s|d')} \quad (6)$$

The policy $\pi_\phi(s | d)$ learns to suggest actions a that optimize the return *while* generating a trajectory achieving descriptor d . Consequently, the descriptor-conditioned actor can exhibit a wide range of descriptors, effectively distilling some of the capabilities of the archive into a single versatile policy.

4.3 Actor-Critic Training

Algorithm 2 Descriptor-conditioned Actor-Critic Training

```

function TRAIN_ACTOR_CRITIC( $\pi_\phi, Q_{\theta_1}, Q_{\theta_2}, \mathcal{B}$ )
  for  $t = 1 \rightarrow n$  do
    Sample  $N$  transitions  $(s, a, r, s', d, d')$  from  $\mathcal{B}$ 
    Sample smoothing noise  $\epsilon$ 
     $y \leftarrow S(d, d') r + \gamma \min_{i=1,2} Q_{\theta_i}(s', \pi_{\phi'}(s' | d') + \epsilon | d')$ 
    Update both critics by regression to  $y$ 
    if  $t \bmod \Delta$  then
      Update actor using the deterministic policy gradient:
       $\frac{1}{N} \sum \nabla_\phi \pi_\phi(s | d') \nabla_a Q_{\theta_1}(s, a | d')|_{a=\pi_\phi(s|d')}$ 
      Soft-update target networks  $Q_{\theta_{i'}}$  and  $\pi_{\phi'}$ 

```

In Section 4.1, we show that the descriptor-conditioned critic target y in Equation (5) requires a transition (s, a, r, s', d) and a target descriptor d' . Most related methods that are conditioned on skills or goals rely on a sampling strategy. For example, HER [1] is a goal-conditioned reinforcement learning algorithm that relies on a handcrafted goal sampling strategy and DIAYN, DADS, SMERL sample skills from a uniform prior distribution. However, in this work, we don't need to rely on an explicit descriptor sampling strategy.

For each PG variation operator offspring, the transitions coming from the evaluation step, are populated with d' equal to the descriptor of the parent solution d_ψ . The PG variation operator mutates the parent to improve fitness while achieving descriptor d_ψ . Thus, although the offspring is not descriptor-conditioned, its implicit target descriptor is d_ψ . Consequently, we set the target descriptor d' to the descriptor of the parent d_ψ .

Similarly, for each GA variation operator offspring, the transitions coming from the evaluation step, are populated with d' equal to the observed descriptor of the trajectory d . The GA variation operator mutates the parent by adding random noise to the genotype. However, a small random change in the parameters of the parent solution can induce big changes in the behavior of the offspring, making them behaviorally different. Consequently, we set the target descriptor d' to the observed descriptor of the trajectory d .

At the end of the evaluation step, we augment the transitions with the observed descriptor of the trajectory d , and with the target descriptor d' , using the implicit descriptor sampling strategy explained above, giving (s, a, r, s', d, d') .

This implicit descriptor sampling strategy has two benefits. First, half of the transitions have $d = d'$, providing the actor-critic training with samples where the target descriptor is achieved, therefore alleviating sparse reward problems. Second, at the beginning of the training process, half of the transitions will have $d \neq d'$ because the solutions in the archive have not learned to accurately achieve their descriptors yet. However, as training goes on, the number of samples where the descriptor is not achieved will decrease, providing some kind of automatic curriculum. Finally, the actor-critic training is adapted from TD3 and is given in Algorithm 2.

4.4 Descriptor-Conditioned PG Variation

Algorithm 3 Descriptor-conditioned PG Variation

```

function VARIATION_PG( $\pi_\psi \dots, Q_{\theta_1}, \mathcal{B}$ )
  for  $\pi_\psi \dots$  do
     $d_\psi \leftarrow D(\pi_\psi)$ 
    for  $i = 1 \rightarrow m$  do
      Sample  $N$  transitions  $(s, a, r, s', d, d')$  from  $\mathcal{B}$ 
      Update actor using the deterministic policy gradient:
       $\frac{1}{N} \sum \nabla_\psi \pi_\psi(s) \nabla_a Q_{\theta_1}(s, a | d_\psi)|_{a=\pi_\psi(s)}$ 
  return  $\pi_{\hat{\phi}} \dots$ 

```

Once the critic $Q_\theta(s, a | d)$ is trained, it can be used to improve the fitness of any solutions in the archive, as described in Algorithm 3. First, a parent solution π_ψ is selected from the archive and we denote its descriptor by $d_\psi := D(\pi_\psi)$. Notice that this policy $\pi_\psi(s)$ is not descriptor-conditioned, contrary to the actor $\pi_\phi(s | d)$. Second, we apply the PG variation operator from Equation (7), for m gradient steps, using the descriptor d_ψ to condition the critic:

$$\nabla_\psi J(\psi) = \frac{1}{N} \sum \nabla_\psi \pi_\psi(s) \nabla_a Q_{\theta_1}(s, a | d_\psi)|_{a=\pi_\psi(s)} \quad (7)$$

The goal is to improve the quality of the solution π_ψ , while keeping the same diversity d_ψ . To that end, the critic is used to evaluate actions and guides π_ψ to (1) improve fitness, while (2) achieving descriptor d_ψ .

4.5 Descriptor-Conditioned Actor Injection

Algorithm 4 Descriptor-conditioned Actor Injection

```

function ACTOR_INJECTION( $\pi_\phi$ )
   $d_1, \dots, d_{b_{AI}} \sim \mathcal{U}(\mathcal{D})$ 
   $\psi_1, \dots, \psi_{b_{AI}} \leftarrow \text{PARAMETERS\_RECOMPUTATION}(\pi_\phi(\cdot | d_1), \dots, \pi_\phi(\cdot | d_{b_{AI}}))$ 
  return  $\pi_{\psi_1}, \dots, \pi_{\psi_{b_{AI}}}$ 

```

In PGA-MAP-ELITES, the actor is injected in the offsprings and considered for addition in the archive at each generation. Empirical analyses [13] have demonstrated the importance of actor injection to achieve good performance. Similarly to PGA-MAP-ELITES, we devise a descriptor-conditioned actor injection (AI) mechanism, to improve the performance of our method, DCG-MAP-ELITES-AI.

There is however a significant challenge. The GA isoline variation operator [44] used in PGA-MAP-ELITES and DCG-MAP-ELITES GECCO requires that all policies in the archive share the same architecture. However, in DCG-MAP-ELITES-AI, the actor is descriptor-conditioned, while the policies in the archive are not. Thus, the first layer of the actor

469 is larger because it takes as input a state and a descriptor, while the first layer of the policies in the archive are smaller
 470 because they take as input only a state. Specifically, for the first layer of the policies in the archive, the weights are a
 471 matrix of dimension $(\dim(\mathcal{S}), 128)$ and the biases are a vector of dimension 128. In contrast, for the first layer of the
 472 descriptor-conditioned actor, the weights are a matrix of dimension $(\dim(\mathcal{S}) + \dim(\mathcal{D}), 128)$ and the biases are a vector
 473 of dimension 128. In both cases, the first hidden layer has 128 neurons, and the subsequent layers are the same.

475 However, for a given fixed descriptor d , we can consider that the constant descriptor d , in $\pi_\phi(s | d)$ is not part
 476 of the input, but part of the parameters. As a matter of fact, for a static descriptor d , we can obtain an equivalent
 477 specialized policy $\pi_{\psi_d}(s)$ with new parameters ψ_d , that is identical to the descriptor-conditioned actor $\pi_\phi(s | d)$, in
 478 terms of state-action mapping. In the following, we show that, given a descriptor d , we can ‘specialize’ the versatile
 479 descriptor-conditioned actor into a non-descriptor-conditioned policy with the same architecture as the policies stored
 480 in the archive. By sampling multiple descriptors, we can perform several actor injections and attempt to add specialized
 481 versions of the versatile actor in niches where it is high-performing, circumventing the need for expensive PG variations.
 482

483 We denote the concatenation operator between two vectors by $||$, the weights and biases of the first layer of the
 484 descriptor-conditioned actor by \mathbf{W} and \mathbf{b} . Given any states s and a descriptor d , we can compute the first layer of the
 485 descriptor-conditioned actor as $(s||d)^\top \mathbf{W} + \mathbf{b} = s^\top \mathbf{W}_1 + (d^\top \mathbf{W}_2 + \mathbf{b})$, with \mathbf{W}_1 a matrix of dimension $(\dim(\mathcal{S}), 128)$
 486 and \mathbf{W}_2 a matrix of dimension $(\dim(\mathcal{D}), 128)$. Therefore, we can reinterpret the computation of the first layer as
 487 the state s multiplied with the matrix \mathbf{W}_1 plus the bias $d^\top \mathbf{W}_2 + \mathbf{b}$. Notice that the matrix \mathbf{W}_1 and bias $d^\top \mathbf{W}_2 + \mathbf{b}$
 488 have the same dimension as the policies in the archive. Thus, if the remaining layers have the same size, we can
 489 recompute the parameters of the first layer, in order to match the architectures and inject the specialized versions of
 490 the descriptor-conditioned actor in the archive.
 491

492 In DCG-MAP-ELITES-AI implementation, we uniformly sample $b_{AI} = 64$ descriptors $d_1, \dots, d_{b_{AI}}$ in the descriptor
 493 space \mathcal{D} . Then, we specialize the descriptor-conditioned actor by recomputing its parameter for each sample descriptor.
 494 At each generation, the resulting policies are suggested for addition in the archive, see Algorithm 4.
 495

496 5 EXPERIMENTS

497 Each experiment is replicated 20 times with random seeds, over one million evaluations and the implementations
 498 are based on the QDax library [3]. The full source code will be made available upon acceptance, in a containerized
 499 environment in which all the experiments and figures can be reproduced. For the quantitative results, we report p-values
 500 based on the Wilcoxon–Mann–Whitney U test with Holm-Bonferroni correction.
 501


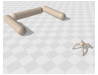





502 5.1 Tasks

503 We evaluate DCG-MAP-ELITES-AI on seven continuous control locomotion QD tasks [31] implemented in Brax [18] and
 504 derived from standard RL benchmarks, see Table 1. Ant Omni, AntTrap Omni and Humanoid Omni are *omnidirectional*
 505 tasks, in which the objective is to minimize energy consumption and the descriptor is the final position of the agent.
 506 Walker Uni, HalfCheetah Uni, Ant Uni and Humanoid Uni are *unidirectional* task in which the objective is to go forward
 507 as fast as possible while minimizing energy consumption and the descriptor is the feet contact rate for each foot of the
 508 agent. Walker Uni, HalfCheetah Uni, Ant Uni were introduced in PGA-MAP-ELITES paper [31] and Humanoid Uni,
 509 Ant Omni, Humanoid Omni were introduced by Flageat et al. [15]. AntTrap Omni is adapted from QD-PG paper [34],
 510 the only difference being the elimination of the forward term in the reward function. We introduce AntTrap Omni to
 511 evaluate DCG-MAP-ELITES-AI on a deceptive, omnidirectional environment. The trap creates a discontinuity of fitness
 512 in the descriptor space as points on both sides of the trap are close, but require two different trajectories to achieve
 513

these descriptors. Thus, the descriptor-conditioned critic needs to learn that discontinuity to provide accurate policy gradients.

PGA-MAP-ELITES has previously shown state-of-the-art results on unidirectional tasks, in particular Walker Uni, HalfCheetah Uni and Ant Uni, but tends to struggle on omnidirectional tasks. In omnidirectional tasks, the global maximum of the fitness function is a solution that does not move, which is directly opposed to discovering how to reach different locations. Hence, the offsprings generated by the PG variation operator will tend to move less and travel a shorter distance. Instead, DCG-MAP-ELITES-AI aims to improve the energy consumption while maintaining the ability to reach distant locations.

Table 1. Evaluation Tasks

	ANT OMNI	ANTTRAP OMNI	HUMANOID OMNI	WALKER UNI	HALFCHEETAH UNI	ANT UNI	HUMANOID UNI
							
STATE	Position and velocity of body parts						
ACTION	Torques applied at the hinge joints						
STATE DIM	30	30	245	18	19	30	245
ACTION DIM	8	8	17	6	6	8	17
DESCRIPTOR DIM	2	2	2	2	2	4	2
EPISODE LEN	250	250	1000	1000	1000	1000	1000
PARAMETERS	21,512	21,512	50,193	19,718	19,846	21,512	50,193

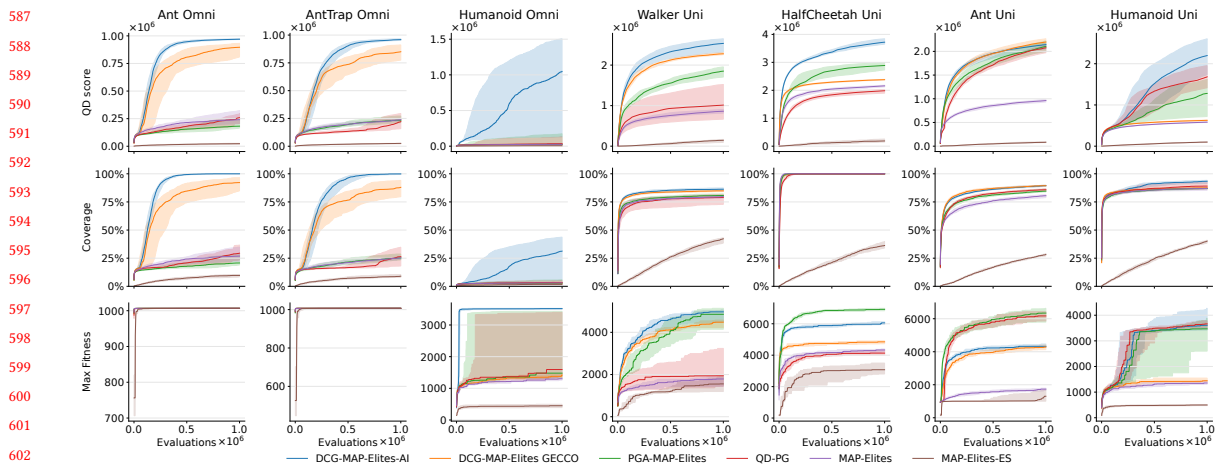
5.2 Main Results

5.2.1 Baselines. We compare DCG-MAP-ELITES-AI with four state-of-the-art algorithms, namely MAP-ELITES [43], MAP-ELITES-ES [7], PGA-MAP-ELITES [31] and QD-PG [34].

5.2.2 Metrics. We consider the QD score, coverage and max fitness to evaluate the final populations (i.e. archives) of all algorithms throughout training, as defined in Flageat et al. [15], Pugh et al. [35] and used in PGA-MAP-ELITES paper [31]. The main metric is the *QD score*, which represents the sum of fitness of all solutions stored in the archive. This metric captures both the quality and the diversity of the population. In the tasks considered, the fitness is always positive, which avoids penalizing algorithms for finding additional solutions. We also consider the *coverage*, which represents the proportion of filled cells in the archive, measuring descriptor space illumination. Finally, we also report the *max fitness*, which is defined as the fitness of the best solution in the archive.

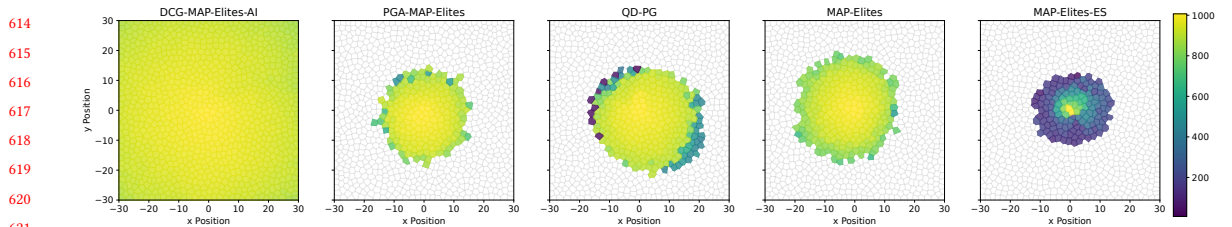
5.2.3 Results. The experimental results presented in Figure 2 demonstrate that DCG-MAP-ELITES-AI achieves equal or higher QD score and coverage than all baselines on all tasks, especially PGA-MAP-ELITES, the previous state-of-the-art. On Ant Uni and Humanoid Uni, DCG-MAP-ELITES-AI achieves a higher median QD score but not significantly. On all other tasks, DCG-MAP-ELITES-AI achieves a significantly higher QD score ($p < 0.003$), demonstrating that our method generates populations of solutions that are higher-performing and more diverse. Especially, the coverage metric shows that DCG-MAP-ELITES-AI surpasses the exploration capabilities of QD-PG on all tasks ($p < 0.05$). DCG-MAP-ELITES-AI significantly outperforms the GECCO version [12] on all environments except Ant Uni ($p < 0.01$), where they perform

573 similarly, showing that the improvements made to the algorithm are beneficial. DCG-MAP-ELITES-AI also achieves equal
 574 or significantly better max fitness on all environments except on HalfCheetah Uni and Ant Uni, where PGA-MAP-ELITES
 575 is better, showing room for improvement. Finally, we also show that our method still benefits from the exploration
 576 power of the GA operator even in deceptive environment like AntTrap Omni. The experimental results confirm that
 577 DCG-MAP-ELITES-AI is able to overcome the limits of PGA-MAP-ELITES on omnidirectional tasks while performing
 578 better on the unidirectional tasks ($p < 0.005$) except Ant Uni where our method is not significantly better. Thus,
 579 confirming the interest of having a descriptor-conditioned gradient to make the PG variation operator fruitful in a wider
 580 range of tasks. Overall, DCG-MAP-ELITES-AI shows competitive performance on all metrics and tasks, hence proving
 581 to be the first successful effort in the QD-RL literature to achieve well on both the unidirectional and omnidirectional
 582 tasks. Previous efforts were usually adapted to either one or the other [31, 34, 42].
 583
 584
 585



604 Fig. 2. QD score, coverage and max fitness (Section 5.2.2) for DCG-MAP-ELITES-AI and all baselines on all tasks. Each experiment is
 605 replicated 20 times with random seeds. The solid line is the median and the shaded area represents the first and third quartiles.
 606

607
 608 Qualitative results in Figure 3 also show that DCG-MAP-ELITES-AI discovers solutions that are more diverse and
 609 higher-performing than other baselines on Ant Omni task. The final archives for all algorithms and on all tasks are
 610 provided in Appendix A.1.
 611



622 Fig. 3. **Ant Omni** Archive at the end of training for all algorithms.
 623
 624

5.3 Ablations

5.3.1 *Ablation studies.* We also compare DCG-MAP-ELITES-AI with three ablations, namely DCG-MAP-ELITES GECCO [12], Ablation AI and Ablation Actor. In DCG-MAP-ELITES GECCO, there is no actor injection, but we perform actor evaluation instead to provide on-policy samples to the TD3 algorithm. In Ablation AI, there is no actor injection and no actor evaluation. In Ablation Actor, the actor is not descriptor-conditioned, removing the archive distillation component, but the critic is still descriptor-conditioned.

5.3.2 *Results.* We perform two ablation experiments to show the importance of actor injection and of the descriptor-conditioned actor. AI proves significantly beneficial in terms of QD score, on all tasks ($p < 0.05$) except Ant Uni where they perform comparably. Having a descriptor-conditioned actor $\pi_{\phi}(\cdot | d)$ rather than a normal actor $\pi_{\phi}(\cdot)$ proves significantly beneficial in terms of QD score, on all tasks ($p < 10^{-4}$), demonstrating that the descriptor-conditioned actor enables archive distillation while being beneficial for the critic’s training. DCG-MAP-ELITES GECCO achieves equal or higher QD score than the AI ablation, showing the importance of on-policy samples. Overall, DCG-MAP-ELITES-AI shows competitive performance on all metrics and tasks compared to the ablations, hence proving the importance of the different enhancements compared to PGA-MAP-ELITES.

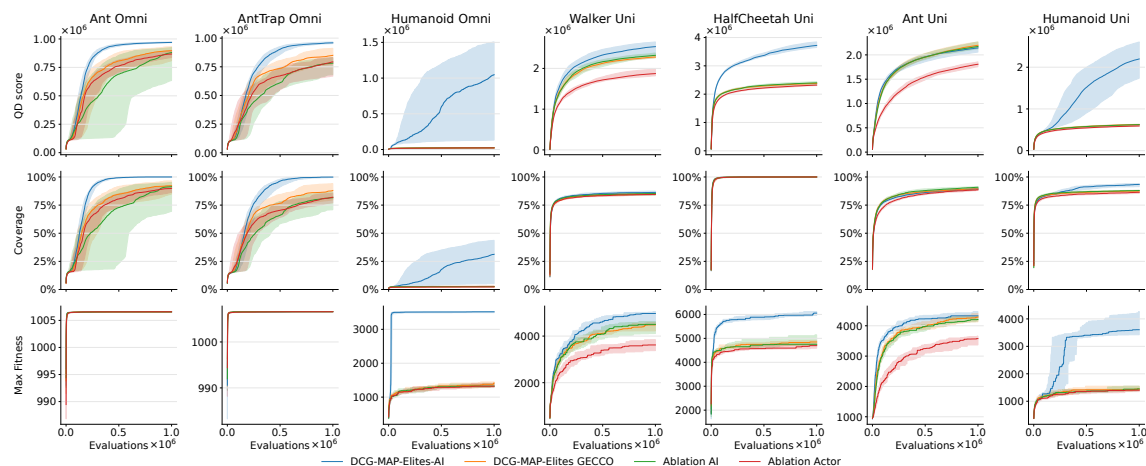


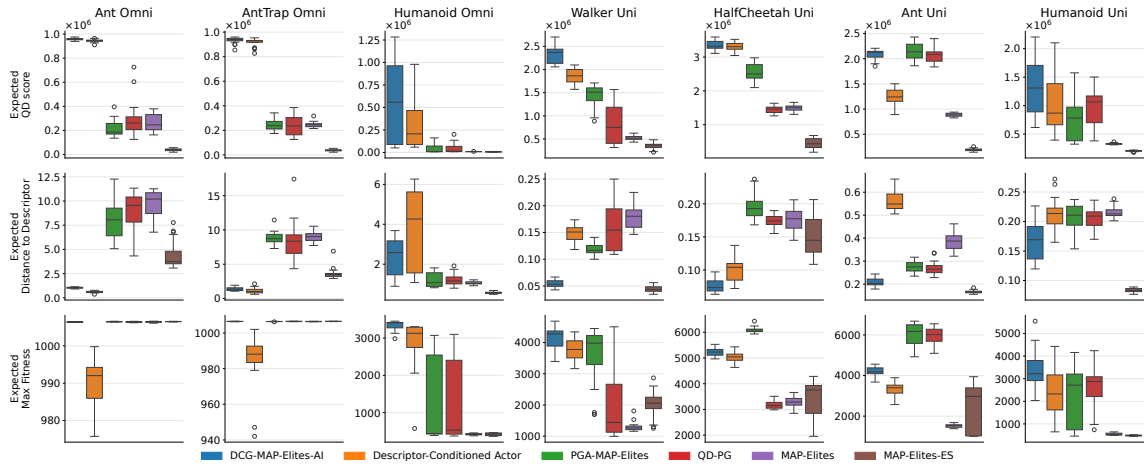
Fig. 4. QD score, coverage and max fitness (Section 5.2.2) for DCG-MAP-ELITES-AI and the ablations on all tasks. Each experiment is replicated 20 times with random seeds. The solid line is the median and the shaded area represents the first and third quartiles.

5.4 Reproducibility

5.4.1 *Reproducibility Metrics.* We also consider three metrics to evaluate the reproducibility of the final archives for all algorithms and of the descriptor-conditioned actor for DCG-MAP-ELITES-AI, at the end of training. QD algorithms based on MAP-ELITES output a population of solutions that we evaluate with the QD score, coverage and max fitness, see Section 5.2.2. However, these metrics can be misleading because in stochastic environments, a solution might give different fitnesses and descriptors when evaluated multiple times. Consequently, the QD score, coverage and max fitness can be overestimated, an effect that is well-known and that has been studied in the past [14]. An archive of solutions is considered reproducible, if the QD score, coverage and max fitness does not change substantially after multiple reevaluation of the individuals. Thus, to assess the reproducibility of the archives, we consider the *expected*

677 *QD score*, the *expected distance to descriptor* and the *expected max fitness*. To calculate those metrics, we reevaluate
 678 each solution in the archive 512 times, to approximate its expected fitness and expected distance to descriptor. The
 679 expected distance to descriptor of a solution is simply the expected euclidean distance between the descriptor of the
 680 cell of the individual and the observed descriptors. Therefore, for the expected distance to descriptor, lower is better.
 681 We use the expected fitness and expected distance to descriptor of all solutions to calculate the expected QD score,
 682 expected distance to descriptor and expected max fitness of the archive.

684 Additionally, DCG-MAP-ELITES-AI's descriptor-conditioned actor can in principle achieve different descriptors and
 685 thus, is comparable to an archive. Similarly to the archive, we evaluate its expected QD score, expected distance to
 686 descriptor and expected max fitness. To that end, we take the descriptor d of each filled cell in the corresponding archive,
 687 and evaluate the actor $\pi_\phi(\cdot | d)$ 512 times, to approximate its expected fitness and expected distance to descriptor.
 688 Analogously to the archive, we use the expected fitnesses and expected distances to descriptor to calculate the expected
 689 QD score, expected distance to descriptor and expected max fitness of the descriptor-conditioned actor.
 690
 691
 692
 693



704
 705
 706
 707
 708
 709
 710
 711 Fig. 5. Expected QD score, expected distance to descriptor (lower is better) and expected max fitness (Section 5.4.1) for DCG-MAP-
 712 ELITES-AI, the descriptor-conditioned policy and the baselines on all tasks. Each experiment is replicated 20 times with random seeds.

713
 714
 715 **5.4.2 Results.** In Figure 5, we provide the expected QD score, expected distance to descriptor and expected max fitness
 716 of the final archive and the descriptor-conditioned policy, see Section 5.4.1. First, we can see that DCG-MAP-ELITES-AI's
 717 final archive achieves equal or higher expected QD score than all baselines on all tasks. The descriptor-conditioned
 718 actor performs similarly to DCG-MAP-ELITES-AI on most environments, but performs significantly worse on Ant
 719 Uni. This shows that, in most cases, the descriptor-conditioned actor is able to restore the quality of the archive
 720 although having compressed the information in a single network. Second, DCG-MAP-ELITES-AI obtains better expected
 721 distance to descriptor (lower is better) than all baselines except MAP-ELITES-ES on all tasks. However, MAP-ELITES-ES
 722 obtains worse QD score and most importantly, worst coverage, making it easier for MAP-ELITES-ES to achieve a low
 723 expected distance to descriptor. DCG-MAP-ELITES-AI descriptor-conditioned actor obtains similar expected distance to
 724 descriptor on omnidirectional. However, it performs consistently worse on unidirectional tasks. This shows that in
 725 some cases, while compressing the quality of the archive in a single network, the descriptor-conditioned actor can
 726
 727
 728

also exhibit the same diversity as the population. Those two combined observations show that the final archive and descriptor-conditioned policy have similar properties on omnidirectional tasks. Overall, those results show that our single descriptor-conditioned policy can already be seen as a promising summary of our archive, showing very similar properties on half our tasks.

5.5 Variation Operators Evaluation

5.5.1 Variation Operator Metrics. DCG-MAP-ELITES-AI and PGA-MAP-ELITES make use of a GA variation operator and of a PG variation operator. The GA variation operator is strictly the same in both algorithms. However, DCG-MAP-ELITES-AI enhances PGA-MAP-ELITES's PG variation operator with a descriptor-conditioned critic, as explained in Section 4.4. To evaluate the performance of each variation operator, we introduce a metric defined as the accumulated number of offsprings added to the archive coming from each variation operator throughout training, that we call *number of elites*. By tracking the number of elites generated by each variation operator over the course of training, we can analyze the interaction and dynamics between the different variation operators and actor injection, providing insights into the relative contributions of the different components.

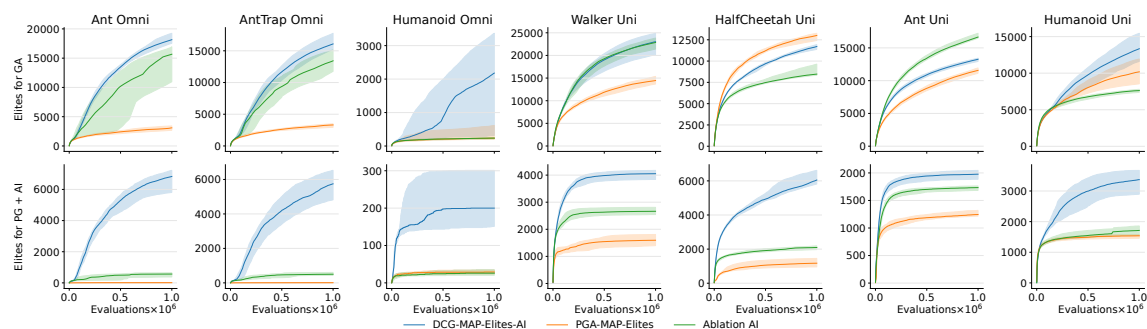


Fig. 6. Accumulated number of offsprings added to the archive (Section 5.5.1) for **(top)** GA variation operator and **(bottom)** PG variation operator plus Actor Injection (AI). Each experiment is replicated 20 times with random seeds. The solid line is the median and the shaded area represents the first and third quartiles.

5.5.2 Results. On the top row of Figure 6, we can see the accumulated number of elites for the GA variation operator for DCG-MAP-ELITES-AI, PGA-MAP-ELITES and ablation AI throughout training. In all three cases, the number of offsprings suggested for addition in the archive is 128. On the bottom row of Figure 6, we can see the accumulated number of elites for the PG variation operator. In all three cases, the number of offsprings suggested for addition in the archive is 128, but for DCG-MAP-ELITES-AI, the PG variation is divided into 64 coming from the actor injection (Section 4.5) and 64 coming from the PG update using the descriptor-conditioned critic (Section 4.4). First, we can see that the ablation of the actor injection generates a larger number of elites than PGA-MAP-ELITES, demonstrating that the descriptor-conditioned critic generates higher-performing and more diverse solution than the traditional critic used in PGA-MAP-ELITES. Furthermore, we can see that DCG-MAP-ELITES-AI with actor injection mechanism generates even more elites than the descriptor-conditioned PG variation operator alone. Interestingly, we can see that the number of elites generated by DCG-MAP-ELITES-AI is higher than PGA-MAP-ELITES, even though the GA variation operators are exactly the same. This demonstrates that the solutions found by the descriptor-conditioned PG variation operator are better stepping stones.

6 CONCLUSION

In this work, we introduce DCG-MAP-ELITES-AI and demonstrate the benefits of having descriptor-conditioned gradients to evolve populations of large neural networks. We concurrently train a descriptor-conditioned actor, as a by-product of the critic’s training, that can achieve a diversity of high-performing behaviors. In turn, we inject the trained descriptor-conditioned actor in the population, despite network architecture differences, speeding-up training even more. Our method, DCG-MAP-ELITES-AI, achieves equal or better performance than all baselines on seven continuous control locomotion tasks. We also show that the synergy between the fitness improvement capabilities of the PG variations and the exploration capabilities of the GA variations is preserved, even in deceptive environments. The descriptor-conditioned actor demonstrates performance that are similar to the discrete archive, summarizing its capabilities into one single neural network and acting as a continuous archive. We think that distilling the archive into a single policy is a promising method as it enables to have less redundancy compared to a discrete archive in which most of the solutions can be similar, especially between close cells. The descriptor-conditioned policy can also negate the burden of dealing with an archive of thousands of solutions in practical applications.

The benefits of combining RL methods with PGA-MAP-ELITES come with the limitations of MDP settings. Specifically, we are limited to evolving differentiable solutions and the foundations of RL algorithms rely on the Markov property and full observability. In this work in particular, we face challenges with the Markov property because the descriptors depend on full trajectories. Thus, the scaled reward introduced in our method depends on the full trajectory and not only on the current state and action. The performance of the descriptor-conditioned policy also shows that there is room for improvement to better distill the knowledge of the archive.

For future work, we would like to investigate the generalization capabilities of the descriptor-conditioned policy trained with DCG-MAP-ELITES-AI and try to produce solutions with descriptors that are not in the archive, performing descriptor space generalization. In our method, the critic attempts to mutate solutions to produce offspring with higher fitness while keeping their descriptors constant. We think that we could use the descriptor-conditioned critic to mutate solutions to produce offspring towards different descriptors, thereby explicitly promoting diversity.

REFERENCES

- [1] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. 2017. Hindsight Experience Replay. In *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/hash/453fadbd8a1a3af50a9df4df899537b5-Abstract.html
- [2] Felix Chalumeau, Raphael Boige, Bryan Lim, Valentin Macé, Maxime Allard, Arthur Flajolet, Antoine Cully, and Thomas Pierrot. 2022. Neuroevolution is a Competitive Alternative to Reinforcement Learning for Skill Discovery. <https://openreview.net/forum?id=6BHLZgyPOZY>
- [3] Felix Chalumeau, Bryan Lim, Raphael Boige, Maxime Allard, Luca Grillotti, Manon Flageat, Valentin Macé, Arthur Flajolet, Thomas Pierrot, and Antoine Cully. 2023. QDax: A Library for Quality-Diversity and Population-based Algorithms with Hardware Acceleration. arXiv:2308.03665 [cs.AI]
- [4] Felix Chalumeau, Thomas Pierrot, Valentin Macé, Arthur Flajolet, Karim Beguir, Antoine Cully, and Nicolas Perrin-Gilbert. 2023. Assessing Quality-Diversity Neuro-Evolution Algorithms Performance in Hard Exploration Problems. <https://doi.org/10.48550/arXiv.2211.13742> arXiv:2211.13742 [cs].
- [5] Konstantinos Chatzilygeroudis, Antoine Cully, Vassilis Vassiliades, and Jean-Baptiste Mouret. 2021. Quality-Diversity Optimization: A Novel Branch of Stochastic Optimization. In *Black Box Optimization, Machine Learning, and No-Free Lunch Theorems*, Panos M. Pardalos, Varvara Rasskazova, and Michael N. Vrahatis (Eds.). Springer International Publishing, Cham, 109–135. https://doi.org/10.1007/978-3-030-66515-9_4
- [6] Konstantinos Chatzilygeroudis, Vassilis Vassiliades, and Jean-Baptiste Mouret. 2018. Reset-free Trial-and-Error Learning for Robot Damage Recovery. *Robotics and Autonomous Systems* 100 (Feb. 2018), 236–250. <https://doi.org/10.1016/j.robot.2017.11.010>
- [7] Cédric Colas, Vashisht Madhavan, Joost Huizinga, and Jeff Clune. 2020. Scaling MAP-Elites to deep neuroevolution. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference (GECCO '20)*. Association for Computing Machinery, New York, NY, USA, 67–75. <https://doi.org/10.1145/3377930.3390217>
- [8] Antoine Cully, Jeff Clune, Danesh Tarapore, and Jean-Baptiste Mouret. 2015. Robots that can adapt like animals. *Nature* 521, 7553 (May 2015), 503–507. <https://doi.org/10.1038/nature14422> Number: 7553 Publisher: Nature Publishing Group.

- 833 [9] Antoine Cully and Yiannis Demiris. 2018. Quality and Diversity Optimization: A Unifying Modular Framework. *IEEE Transactions on Evolutionary*
834 *Computation* 22, 2 (2018), 245–259. <https://doi.org/10.1109/TEVC.2017.2704781>
- 835 [10] Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O. Stanley, and Jeff Clune. 2021. First return, then explore. *Nature* 590, 7847 (Feb. 2021),
836 580–586. <https://doi.org/10.1038/s41586-020-03157-9> Number: 7847 Publisher: Nature Publishing Group.
- 837 [11] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. 2018. Diversity is All You Need: Learning Skills without a Reward Function.
838 <https://doi.org/10.48550/arXiv.1802.06070> arXiv:1802.06070 [cs].
- 839 [12] Maxence Faldor, Félix Chalumeau, Manon Flageat, and Antoine Cully. 2023. MAP-Elites with Descriptor-Conditioned Gradients and Archive
840 Distillation into a Single Policy. In *Proceedings of the Genetic and Evolutionary Computation Conference (Lisbon, Portugal) (GECCO '23)*. Association
841 for Computing Machinery, New York, NY, USA, 138–146. <https://doi.org/10.1145/3583131.3590503>
- 842 [13] Manon Flageat, Félix Chalumeau, and Antoine Cully. 2023. Empirical analysis of PGA-MAP-Elites for Neuroevolution in Uncertain Domains. *ACM*
843 *Transactions on Evolutionary Learning and Optimization* 3, 1 (March 2023), 1:1–1:32. <https://doi.org/10.1145/3577203>
- 844 [14] Manon Flageat and Antoine Cully. 2023. Uncertain Quality-Diversity: Evaluation methodology and new methods for Quality-Diversity in Uncertain
845 Domains. *IEEE Transactions on Evolutionary Computation* (2023).
- 846 [15] Manon Flageat, Bryan Lim, Luca Grillotti, Maxime Allard, Simón C. Smith, and Antoine Cully. 2022. Benchmarking Quality-Diversity Algorithms
847 on Neuroevolution for Reinforcement Learning. <https://doi.org/10.48550/arXiv.2211.02193> arXiv:2211.02193 [cs].
- 848 [16] Matthew Fontaine and Stefanos Nikolaidis. 2021. Differentiable Quality Diversity. In *Advances in Neural Information Processing Systems*, Vol. 34.
849 Curran Associates, Inc., 10040–10052. <https://proceedings.neurips.cc/paper/2021/hash/532923f11ac97d3e7cb0130315b067dc-Abstract.html>
- 850 [17] Matthew Fontaine and Stefanos Nikolaidis. 2023. Covariance Matrix Adaptation MAP-Annealing. In *Proceedings of the Genetic and Evolutionary*
851 *Computation Conference (GECCO '23)*. Association for Computing Machinery, New York, NY, USA, 456–465. <https://doi.org/10.1145/3583131.3590389>
- 852 [18] C. Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. 2021. *Brax - A Differentiable Physics Engine for*
853 *Large Scale Rigid Body Simulation*. <http://github.com/google/brax>
- 854 [19] Scott Fujimoto, Herke Hoof, and David Meger. 2018. Addressing Function Approximation Error in Actor-Critic Methods. In *Proceedings of the 35th*
855 *International Conference on Machine Learning*. PMLR, 1587–1596. <https://proceedings.mlr.press/v80/fujimoto18a.html> ISSN: 2640-3498.
- 856 [20] Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. 2016. Variational Intrinsic Control. <https://doi.org/10.48550/arXiv.1611.07507>
857 arXiv:1611.07507 [cs].
- 858 [21] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement
859 Learning with a Stochastic Actor. In *Proceedings of the 35th International Conference on Machine Learning*. PMLR, 1861–1870. <https://proceedings.mlr.press/v80/haarnoja18b.html> ISSN: 2640-3498.
- 860 [22] Nikolaus Hansen. 2023. The CMA Evolution Strategy: A Tutorial. <https://doi.org/10.48550/arXiv.1604.00772> arXiv:1604.00772 [cs, stat].
- 861 [23] Nicolas Heess, Dhruva TB, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne, Yuval Tassa, Tom Erez, Ziyu Wang, Ali Eslami, Martin
862 Riedmiller, and David Silver. 2017. Emergence of Locomotion Behaviours in Rich Environments. (July 2017).
- 863 [24] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. 1989. Multilayer feedforward networks are universal approximators. *Neural Networks* 2, 5
864 (1989), 359–366. [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8)
- 865 [25] Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M. Czarnecki, Jeff Donahue, Ali Razavi, Oriol Vinyals, Tim Green, Iain Dunning,
866 Karen Simonyan, Chrisantha Fernando, and Koray Kavukcuoglu. 2017. Population Based Training of Neural Networks. [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.1711.09846)
867 arXiv:1711.09846 arXiv:1711.09846 [cs].
- 868 [26] Marija Jegorova, Stéphane Doncieux, and Timothy Hospedales. 2019. Behavioural Repertoire via Generative Adversarial Policy Networks. In
869 *2019 Joint IEEE 9th International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*. [https://doi.org/10.1109/ICDL-](https://doi.org/10.1109/ICDL-EpiRob44920.2019)
870 EpiRob44920.2019 arXiv:1811.02945 [cs, stat].
- 871 [27] Saurabh Kumar, Aviral Kumar, Sergey Levine, and Chelsea Finn. 2020. One Solution is Not All You Need: Few-Shot Extrapolation via Structured
872 MaxEnt RL. In *Advances in Neural Information Processing Systems*, Vol. 33. Curran Associates, Inc., 8198–8210. [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2020/hash/5d151d1059a6281335a10732fc49620e-Abstract.html)
873 2020/hash/5d151d1059a6281335a10732fc49620e-Abstract.html
- 874 [28] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2016. Continuous
875 control with deep reinforcement learning. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4,*
876 *2016, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1509.02971>
- 877 [29] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K.
878 Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra,
879 Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (Feb. 2015), 529–533.
880 <https://doi.org/10.1038/nature14236> Number: 7540 Publisher: Nature Publishing Group.
- 881 [30] Jean-Baptiste Mouret and Jeff Clune. 2015. Illuminating search spaces by mapping elites. *CoRR* abs/1504.04909 (2015). arXiv:1504.04909 <http://arxiv.org/abs/1504.04909>
- 882 [31] Olle Nilsson and Antoine Cully. 2021. Policy gradient assisted MAP-Elites. In *Proceedings of the Genetic and Evolutionary Computation Conference*
883 *(GECCO '21)*. Association for Computing Machinery, New York, NY, USA, 866–875. <https://doi.org/10.1145/3449639.3459304>
- 884 [32] OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn
885 Powell, Raphael Ribas, Jonas Schneider, Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba, and Lei Zhang.
886 2019. Solving Rubik’s Cube with a Robot Hand. <https://doi.org/10.48550/arXiv.1910.07113> arXiv:1910.07113 [cs, stat].

- 885 [33] Thomas Pierrot and Arthur Flajolet. 2023. Evolving Populations of Diverse RL Agents with MAP-Elites. <https://doi.org/10.48550/arXiv.2303.12803>
886 arXiv:2303.12803 [cs].
- 887 [34] Thomas Pierrot, Valentin Macé, Felix Chalumeau, Arthur Flajolet, Geoffrey Cideron, Karim Beguir, Antoine Cully, Olivier Sigaud, and Nicolas Perrin-
888 Gilbert. 2022. Diversity policy gradient for sample efficient quality-diversity optimization. In *Proceedings of the Genetic and Evolutionary Computation*
889 *Conference (GECCO '22)*. Association for Computing Machinery, New York, NY, USA, 1075–1083. <https://doi.org/10.1145/3512290.3528845>
- 890 [35] Justin K. Pugh, Lisa B. Soros, and Kenneth O. Stanley. 2016. Quality Diversity: A New Frontier for Evolutionary Computation. *Frontiers in Robotics*
891 *and AI* 3 (2016). <https://www.frontiersin.org/articles/10.3389/frobt.2016.00040>
- 892 [36] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. 2017. Evolution strategies as a scalable alternative to reinforcement learning.
893 *arXiv preprint arXiv:1703.03864* (2017).
- 894 [37] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. 2015. Universal Value Function Approximators. In *Proceedings of the 32nd International*
895 *Conference on Machine Learning*. PMLR, 1312–1320. <https://proceedings.mlr.press/v37/schaul15.html> ISSN: 1938-7228.
- 896 [38] Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. 2019. Dynamics-Aware Unsupervised Discovery of Skills.
897 <https://openreview.net/forum?id=HJgLR4KvH>
- 898 [39] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda
899 Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine
900 Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. 2016. Mastering the game of Go with deep neural networks and tree search. *Nature*
901 529, 7587 (Jan. 2016), 484–489. <https://doi.org/10.1038/nature16961> Number: 7587 Publisher: Nature Publishing Group.
- 902 [40] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. 2014. Deterministic Policy Gradient Algorithms.
903 In *Proceedings of the 31st International Conference on Machine Learning*. PMLR, 387–395. <https://proceedings.mlr.press/v32/silver14.html> ISSN:
904 1938-7228.
- 905 [41] Bryon Tjanaka, Matthew C. Fontaine, David H. Lee, Aniruddha Kalkar, and Stefanos Nikolaidis. 2023. Training Diverse High-Dimensional Controllers
906 by Scaling Covariance Matrix Adaptation MAP-Annealing. <https://doi.org/10.48550/arXiv.2210.02622> arXiv:2210.02622 [cs].
- 907 [42] Bryon Tjanaka, Matthew C. Fontaine, Julian Togelius, and Stefanos Nikolaidis. 2022. Approximating gradients for differentiable quality diversity in
908 reinforcement learning. In *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO '22)*. Association for Computing Machinery,
909 New York, NY, USA, 1102–1111. <https://doi.org/10.1145/3512290.3528705>
- 910 [43] Vassilis Vassiliades, Konstantinos Chatzilygeroudis, and Jean-Baptiste Mouret. 2018. Using Centroidal Voronoi Tessellations to Scale Up the
911 Multidimensional Archive of Phenotypic Elites Algorithm. *IEEE Transactions on Evolutionary Computation* 22, 4 (2018), 623–630. <https://doi.org/10.1109/TEVC.2017.2735550>
- 912 [44] Vassiliis Vassiliades and Jean-Baptiste Mouret. 2018. Discovering the elite hypervolume by leveraging interspecies correlation. In *Proceedings*
913 *of the Genetic and Evolutionary Computation Conference (GECCO '18)*. Association for Computing Machinery, New York, NY, USA, 149–156.
914 <https://doi.org/10.1145/3205455.3205602>
- 915 [45] Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo
916 Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max
917 Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine,
918 Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul,
919 Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. 2019. Grandmaster level in StarCraft II using multi-agent
920 reinforcement learning. *Nature* 575, 7782 (Nov. 2019), 350–354. <https://doi.org/10.1038/s41586-019-1724-z> Number: 7782 Publisher: Nature
921 Publishing Group.
- 922
- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936

A SUPPLEMENTARY RESULTS

A.1 Archives

We provide the archives obtained at the end of training for each algorithm on all environments. For each (algorithm, environment) pair, we select the most representative seed with the QD score closest to the median QD score over all seeds to avoid cherry picking.

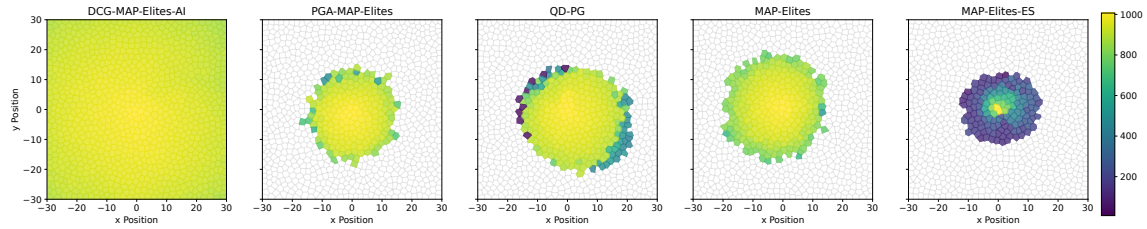


Fig. 7. **Ant Omni** Archive at the end of training for all algorithms.

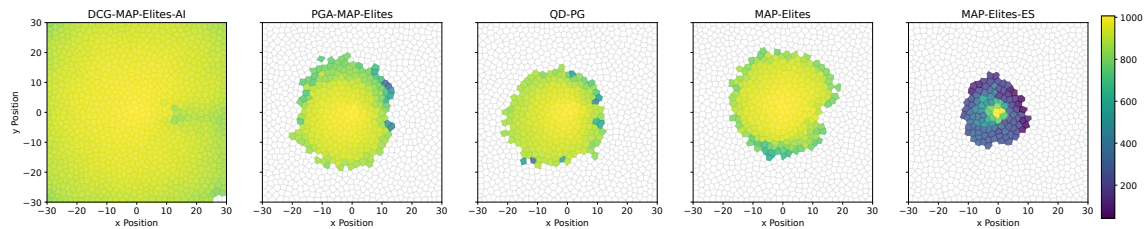


Fig. 8. **AntTrap Omni** Archive at the end of training for all algorithms.

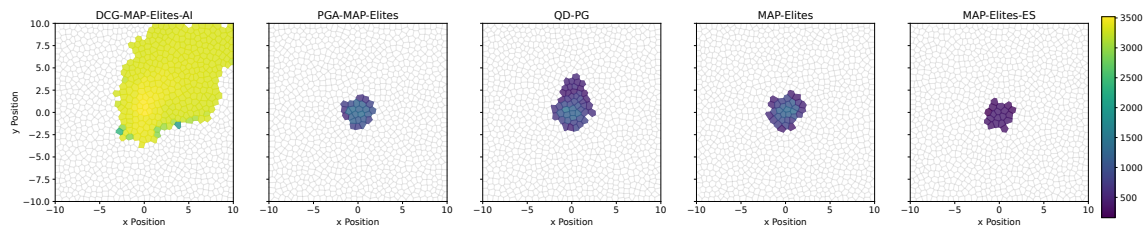


Fig. 9. **Humanoid Omni** Archive at the end of training for all algorithms.

989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040

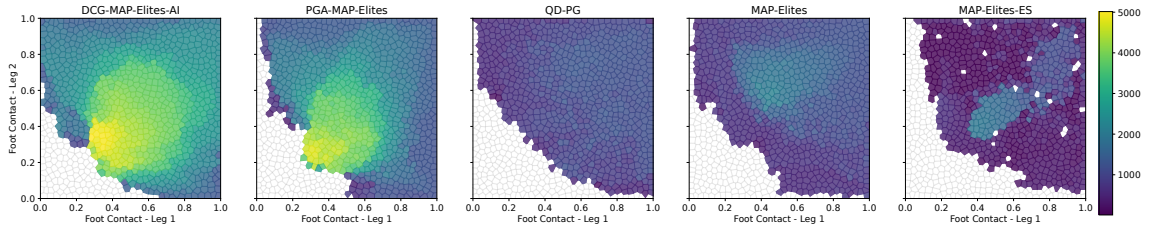


Fig. 10. Walker Uni Archive at the end of training for all algorithms.

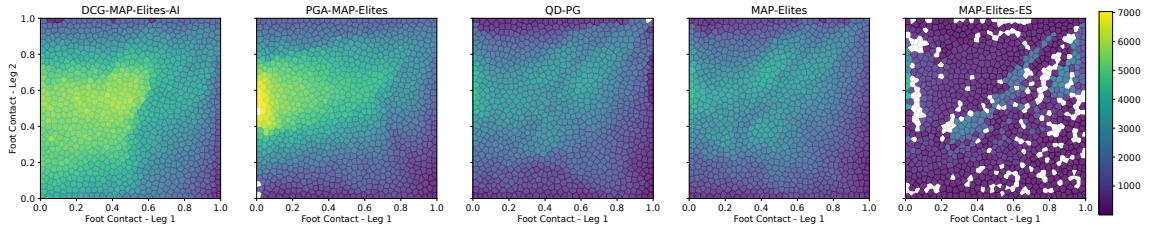


Fig. 11. Halfcheetah Uni Archive at the end of training for all algorithms.

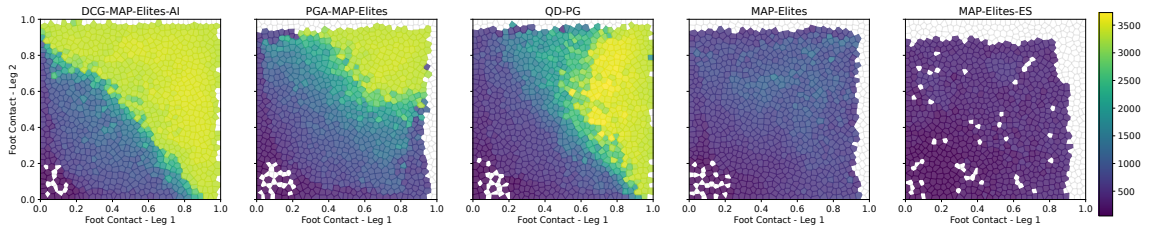


Fig. 12. Humanoid Uni Archive at the end of training for all algorithms.

1041 B ALGORITHMS

1042 B.1 DCG-MAP-ELITES-AI

1044 Algorithm 5 DCG-MAP-ELITES-AI

1046 **Require:** GA batch size b_{GA} , PG batch size b_{PG} , Actor Injection batch size b_{AI} , total batch size $b = b_{GA} + b_{PG} + b_{AI}$

1047 Initialize archive \mathcal{X} with b random solutions and replay buffer \mathcal{B}

1048 Initialize critic networks $Q_{\theta_1}, Q_{\theta_2}$ and actor network π_{ϕ}

1049 $i \leftarrow 0$

1050 **while** $i < I$ **do**

1051 TRAIN_ACTOR_CRITIC($\pi_{\phi}, Q_{\theta_1}, Q_{\theta_2}, \mathcal{B}$)

1052 $\pi_{\psi_1}, \dots, \pi_{\psi_b} \leftarrow \text{SELECTION}(\mathcal{X})$

1053 $\pi_{\hat{\psi}_1}, \dots, \pi_{\hat{\psi}_{b_{GA}}} \leftarrow \text{VARIATION_GA}(\pi_{\psi_1}, \dots, \pi_{\psi_{b_{GA}}})$

1054 $\pi_{\hat{\psi}_{b_{GA}+1}}, \dots, \pi_{\hat{\psi}_{b_{GA}+b_{PG}}} \leftarrow \text{VARIATION_PG}(\pi_{\psi_{b_{GA}+1}}, \dots, \pi_{\psi_{b_{GA}+b_{PG}}}, Q_{\theta_1}, \mathcal{B})$

1055 $\pi_{\hat{\psi}_{b_{GA}+b_{PG}+1}}, \dots, \pi_{\hat{\psi}_b} \leftarrow \text{ACTOR_INJECTION}(\pi_{\phi})$

1056 ADDITION($\pi_{\hat{\psi}_1}, \dots, \pi_{\hat{\psi}_b}, \mathcal{X}, \mathcal{B}$)

1057 $i \leftarrow i + b$

1058 **function** ADDITION($\pi_{\hat{\psi}} \dots, \mathcal{X}, \mathcal{B}$)

1059 **for** $\pi_{\hat{\psi}} \dots$ **do**

1060 $(f, \text{transitions}) \leftarrow F(\pi_{\hat{\psi}}), d \leftarrow D(\pi_{\hat{\psi}})$

1061 INSERT($\mathcal{B}, \text{transitions}$)

1062 **if** $\mathcal{X}(d) = \emptyset$ or $F(\mathcal{X}(d)) < f$ **then**

1063 $\mathcal{X}(d) \leftarrow \pi_{\hat{\psi}}$

1064 Algorithm 6 Descriptor-conditioned Actor-Critic Training

1065 **function** TRAIN_ACTOR_CRITIC($\pi_{\phi}, Q_{\theta_1}, Q_{\theta_2}, \mathcal{B}$)

1066 **for** $t = 1 \rightarrow n$ **do**

1067 Sample N transitions (s, a, r, s', d, d') from \mathcal{B}

1068 Sample smoothing noise ϵ

1069 $y \leftarrow S(d, d') r + \gamma \min_{i=1,2} Q_{\theta_i}(s', \pi_{\phi'}(s' | d') + \epsilon | d')$

1070 Update both critics by regression to y

1071 **if** $t \bmod \Delta$ **then**

1072 Update actor using the deterministic policy gradient:

1073 $\frac{1}{N} \sum \nabla_{\phi} \pi_{\phi}(s | d') \nabla_a Q_{\theta_1}(s, a | d')|_{a=\pi_{\phi}(s|d')}$

1074 Soft-update target networks $Q_{\theta_i'}$ and $\pi_{\phi'}$

Algorithm 7 Descriptor-conditioned PG Variation

```

1093 function VARIATION_PG( $\pi_\psi \dots Q_{\theta_1}, \mathcal{B}$ )
1094   for  $\pi_\psi \dots$  do
1095      $d_\psi \leftarrow D(\pi_\psi)$ 
1096     for  $i = 1 \rightarrow m$  do
1097       Sample  $N$  transitions  $(s, a, r, s', d, d')$  from  $\mathcal{B}$ 
1098       Update actor using the deterministic policy gradient:
1099        $\frac{1}{N} \sum \nabla_{\psi} \pi_{\psi}(s) \nabla_a Q_{\theta_1}(s, a | d_{\psi})|_{a=\pi_{\psi}(s)}$ 
1100
1101   return  $\pi_{\hat{\phi}} \dots$ 

```

Algorithm 8 Descriptor-conditioned Actor Injection

```

1107 function ACTOR_INJECTION( $\pi_\phi$ )
1108    $d_1, \dots, d_{b_{AI}} \sim \mathcal{U}(\mathcal{D})$ 
1109    $\psi_1, \dots, \psi_{b_{AI}} \leftarrow \text{PARAMETERS\_RECOMPUTATION}(\pi_\phi(\cdot | d_1), \dots, \pi_\phi(\cdot | d_{b_{AI}}))$ 
1110   return  $\pi_{\psi_1}, \dots, \pi_{\psi_{b_{AI}}}$ 

```

B.2 PGA-MAP-ELITES**Algorithm 9** PGA-MAP-ELITES

```

1120 Require: GA batch size  $b_{GA}$ , PG batch size  $b_{PG}$ , total batch size  $b = b_{GA} + b_{PG}$ 
1121 Initialize archive  $\mathcal{X}$  with  $b$  random solutions and replay buffer  $\mathcal{B}$ 
1122 Initialize critic networks  $Q_{\theta_1}, Q_{\theta_2}$  and actor network  $\pi_\phi$ 
1123  $i \leftarrow 0$ 
1124 while  $i < I$  do
1125   TRAIN_ACTOR_CRITIC( $\pi_\phi, Q_{\theta_1}, Q_{\theta_2}, \mathcal{B}$ )
1126    $\pi_{\psi_1}, \dots, \pi_{\psi_{b-1}} \leftarrow \text{SELECTION}(\mathcal{X})$ 
1127    $\pi_{\hat{\psi}_1}, \dots, \pi_{\hat{\psi}_{b_{GA}}} \leftarrow \text{VARIATION\_GA}(\pi_{\psi_1}, \dots, \pi_{\psi_{b_{GA}}})$ 
1128    $\pi_{\hat{\psi}_{b_{GA}+1}}, \dots, \pi_{\hat{\psi}_{b-1}} \leftarrow \text{VARIATION\_PG}(\pi_{\psi_{b_{GA}+1}}, \dots, \pi_{\psi_{b-1}}, Q_{\theta_1}, \mathcal{B})$ 
1129    $\pi_{\hat{\psi}_b} \leftarrow \text{ACTOR\_INJECTION}(\pi_\phi)$ 
1130   ADDITION( $\mathcal{X}, \pi_{\hat{\psi}_1}, \dots, \pi_{\hat{\psi}_{b-1}}, \pi_\phi, \mathcal{B}$ )
1131    $i \leftarrow i + b$ 
1132 function ADDITION( $\mathcal{X}, \pi_{\hat{\psi}} \dots, \mathcal{B}$ )
1133   for  $\pi_{\hat{\psi}} \dots$  do
1134      $(f, \text{transitions}) \leftarrow F(\pi_{\hat{\psi}}), d \leftarrow D(\pi_{\hat{\psi}})$ 
1135     INSERT( $\mathcal{B}, \text{transitions}$ )
1136     if  $\mathcal{X}(d) = \emptyset$  or  $F(\mathcal{X}(d)) < f$  then
1137        $\mathcal{X}(d) \leftarrow \pi_{\hat{\psi}}$ 

```

Algorithm 10 Actor-Critic Training

function TRAIN_ACTOR_CRITIC($\pi_\phi, Q_{\theta_1}, Q_{\theta_2}, \mathcal{B}$)**for** $t = 1 \rightarrow n$ **do** Sample N transitions (s, a, r, s') from \mathcal{B} Sample smoothing noise ϵ $y \leftarrow r + \gamma \min_{i=1,2} Q_{\theta_i}(s', \pi_{\phi'}(s')) + \epsilon$ Update both critics by regression to y **if** $t \bmod \Delta$ **then**

Update actor using the deterministic policy gradient:

 $\frac{1}{N} \sum \nabla_\phi \pi_\phi(s) \nabla_a Q_{\theta_1}(s, a)|_{a=\pi_\phi(s)}$ Soft-update target networks $Q_{\theta_i'}$ and $\pi_{\phi'}$

Algorithm 11 PG Variation

function VARIATION_PG($\pi_\psi \dots, Q_{\theta_1}, \mathcal{B}$)**for** $\pi_\psi \dots$ **do** **for** $i = 1 \rightarrow m$ **do** Sample N transitions (s, a, r, s') from \mathcal{B}

Update actor using the deterministic policy gradient:

 $\frac{1}{N} \sum \nabla_\psi \pi_\psi(s) \nabla_a Q_{\theta_1}(s, a)|_{a=\pi_\psi(s)}$ **return** $\pi_{\hat{\psi}} \dots$

Algorithm 12 Actor Injection

function ACTOR_INJECTION(π_ϕ)**return** π_ϕ

1197 **B.3 QD-PG**

1198

1199

Algorithm 13 QD-PG

1200

Require: GA batch size b_{GA} , QPG batch size b_{QPG} , DPG batch size b_{DPG} , total batch size $b = b_{GA} + b_{QPG} + b_{DPG}$

1201

Initialize archive \mathcal{X} with b random solutions and replay buffer \mathcal{B}

1202

Initialize critic networks $Q_{\theta_Q}, Q_{\theta_D}$ and actor network π_ϕ

1203

 $i \leftarrow 0$

1204

while $i < I$ **do**

1205

 TRAIN_ACTOR_CRITIC($\pi_\phi, Q_{\theta_Q}, Q_{\theta_D}, \mathcal{B}$)

1206

 $\pi_{\psi_1}, \dots, \pi_{\psi_b} \leftarrow \text{SELECTION}(\mathcal{X})$

1207

 $\pi_{\hat{\psi}_1}, \dots, \pi_{\hat{\psi}_{b_{GA}}} \leftarrow \text{VARIATION_GA}(\pi_{\psi_1}, \dots, \pi_{\psi_{b_{GA}}})$

1208

 $\pi_{\hat{\psi}_{b_{GA}+1}}, \dots, \pi_{\hat{\psi}_{b_{GA}+b_{QPG}}} \leftarrow \text{VARIATION_QPG}(\pi_{\psi_{b_{GA}+1}}, \dots, \pi_{\psi_{b_{GA}+b_{QPG}}}, Q_{\theta_Q}, \mathcal{B})$

1209

 $\pi_{\hat{\psi}_{b_{GA}+b_{QPG}+1}}, \dots, \pi_{\hat{\psi}_b} \leftarrow \text{VARIATION_DPG}(\pi_{\psi_{b_{GA}+b_{QPG}+1}}, \dots, \pi_{\psi_b}, Q_{\theta_D}, \mathcal{B})$

1210

 ADDITION($\pi_{\hat{\psi}_1}, \dots, \pi_{\hat{\psi}_b}, \mathcal{X}, \mathcal{B}$)

1211

 $i \leftarrow i + b$

1212

function ADDITION($\mathcal{X}, \mathcal{B}, \pi_\phi, \pi_{\hat{\psi}} \dots$)

1213

for $d' \in \mathcal{D}$ sampled from b solutions in \mathcal{X} **do**

1214

 $(f, \text{transitions}) \leftarrow F(\pi_\phi(\cdot | d'))$

1215

 INSERT($\mathcal{B}, \text{transitions}$)

1216

for $\pi_{\hat{\psi}} \dots$ **do**

1217

 $(f, \text{transitions}) \leftarrow F(\pi_{\hat{\psi}}), d \leftarrow D(\pi_{\hat{\psi}})$

1218

 INSERT($\mathcal{B}, \text{transitions}$)

1219

if $\mathcal{X}(d) = \emptyset$ or $F(\mathcal{X}(d)) < f$ **then**

1220

 $\mathcal{X}(d) \leftarrow \pi_{\hat{\psi}}$

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249 **B.4 MAP-ELITES**

1250

1251 **Algorithm 14** MAP-ELITES

1252

1253 **Require:** GA batch size b_{GA} 1254 Initialize archive \mathcal{X} with b_{GA} random solutions1255 $i \leftarrow 0$ 1256 **while** $i < I$ **do**1257 $x_1, \dots, x_{b_{GA}} \leftarrow \text{SELECTION}(\mathcal{X})$ 1258 $\widehat{x}_1, \dots, \widehat{x}_{b_{GA}} \leftarrow \text{VARIATION}(x_1, \dots, x_{b_{GA}})$ 1259 $\text{ADDITION}(\mathcal{X}, \widehat{x}_1, \dots, \widehat{x}_{b_{GA}})$ 1260 $i \leftarrow i + b_{GA}$ 1261 **function** $\text{ADDITION}(\mathcal{X}, \widehat{x} \dots)$:1262 **for** $\widehat{x} \dots$ **do**1263 $f \leftarrow F(\widehat{x}), d \leftarrow D(\widehat{x})$ 1264 **if** $\mathcal{X}(d) = \emptyset$ **or** $F(\mathcal{X}(d)) < f$ **then**1265 $\mathcal{X}(d) \leftarrow \widehat{x}$

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301 **B.5 MAP-ELITES-ES**1302
1303 **Algorithm 15** MAP-ELITES-ES1304 **Require:** Number of ES samples N , standard deviation of ES samples σ , explore-exploit alternation N_{gen} , number of
1305 re-sampling M 1306 Initialize archive \mathcal{X} with N random solutions, initialise empty novelty archive \mathcal{A} 1307 $i \leftarrow 0$ 1308 **while** $i < I$ **do**1309 **if** $i \% N_{gen} == 0$ **then:**1310 $x \leftarrow \text{SELECTION_EXPLOIT}(\mathcal{X})$ 1311 $\hat{x} \leftarrow \text{VARIATION_EXPLOIT}(x)$ 1312 **else:**1313 $x \leftarrow \text{SELECTION_EXPLORE}(\mathcal{X})$ 1314 $\hat{x} \leftarrow \text{VARIATION_EXPLORE}(\mathcal{A}, x)$ 1315 $\text{ADDITION}(\mathcal{X}, \mathcal{A}, \hat{x})$ 1316 $i \leftarrow i + N + M$ 1317 **function** $\text{ADDITION}(\mathcal{X}, \mathcal{A}, \hat{x})$:1318 **for** $i = 1, \dots, M$ **do**1319 $f_i \leftarrow F(\hat{x}), d_i \leftarrow D(\hat{x})$ 1320 $f \leftarrow \text{average}(f_i), d \leftarrow \text{average}(d_i)$ 1321 $\mathcal{A} \leftarrow \mathcal{A} + d$ 1322 **if** $\mathcal{X}(d) = \emptyset$ or $F(\mathcal{X}(d)) < f$ **then**1323 $\mathcal{X}(d) \leftarrow \hat{x}$ 1324 **function** $\text{VARIATION_EXPLOIT}(x)$:1325 $x_1, \dots, x_N \leftarrow \text{SAMPLE_GAUSSIAN}(x, \sigma)$ 1326 $f_1, \dots, f_N \leftarrow F(x_1, \dots, x_N)$ 1327 $\hat{x} \leftarrow \text{ES_STEP}(x, f_1, \dots, f_N)$ 1328 **function** $\text{VARIATION_EXPLORE}(\mathcal{A}, x)$:1329 $x_1, \dots, x_N \leftarrow \text{SAMPLE_GAUSSIAN}(x, \sigma)$ 1330 $d_1, \dots, d_N \leftarrow D(x_1, \dots, x_N)$ 1331 $nov_1, \dots, nov_N \leftarrow \text{NOVELTY}(\mathcal{A}, d_1, \dots, d_N)$ 1332 $\hat{x} \leftarrow \text{ES_STEP}(x, nov_1, \dots, nov_N)$ 1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352

C HYPERPARAMETERS

C.1 DCG-MAP-ELITES-AI

Table 2. DCG-MAP-ELITES-AI hyperparameters

Parameter	Value
Number of centroids	1024
Total batch size b	256
GA batch size b_{GA}	128
PG batch size b_{PG}	64
AI batch size b_{AI}	64
Policy networks	$[128, 128, \mathcal{A}]$
GA variation param. 1 σ_1	0.005
GA variation param. 2 σ_2	0.05
Actor network	$[128, 128, \mathcal{A}]$
Critic network	$[256, 256, 1]$
TD3 batch size N	100
Critic training steps n	3000
PG training steps m	150
Policy learning rate	5×10^{-3}
Actor learning rate	3×10^{-4}
Critic learning rate	3×10^{-4}
Replay buffer size	10^6
Discount factor γ	0.99
Actor delay Δ	2
Target update rate	0.005
Smoothing noise var. σ	0.2
Smoothing noise clip	0.5
Length scale l	0.1

1405 **C.2 PGA-MAP-ELITES**

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

Table 3. PGA-MAP-ELITES hyperparameters

Parameter	Value
Number of centroids	1024
Total batch size b	256
GA batch size b_{GA}	128
PG batch size b_{PG}	127
AI batch size b_{AI}	1
Policy networks	$[128, 128, \mathcal{A}]$
GA variation param. 1 σ_1	0.005
GA variation param. 2 σ_2	0.05
Actor network	$[128, 128, \mathcal{A}]$
Critic network	$[256, 256, 1]$
TD3 batch size N	100
Critic training steps n	3000
PG training steps m	150
Policy learning rate	5×10^{-3}
Actor learning rate	3×10^{-4}
Critic learning rate	3×10^{-4}
Replay buffer size	10^6
Discount factor γ	0.99
Actor delay Δ	2
Target update rate	0.005
Smoothing noise var. σ	0.2
Smoothing noise clip	0.5

1457 **C.3 QD-PG**
 1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508

Table 4. QD-PG hyperparameters

Parameter	Value
Number of centroids	1024
Total batch size b	256
GA batch size b_{GA}	86
QPG batch size b_{PG}	85
DPG batch size b_{PG}	85
Policy networks	$[128, 128, \mathcal{A}]$
GA variation param. 1 σ_1	0.005
GA variation param. 2 σ_2	0.05
Actor network	$[128, 128, \mathcal{A}]$
Critic network	$[256, 256, 1]$
TD3 batch size N	100
Quality critic training steps n	3000
Diversity critic training steps n	300
PG training steps m	150
Policy learning rate	5×10^{-3}
Actor learning rate	3×10^{-4}
Critic learning rate	3×10^{-4}
Replay buffer size	10^6
Discount factor γ	0.99
Actor delay Δ	2
Target update rate	0.005
Smoothing noise var. σ	0.2
Smoothing noise clip	0.5
Number nearest neighbors	5
Novelty scaling ratio	1.0

1509 **C.4 MAP-ELITES**

1510

1511

1512

1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1523

1524

1525

C.5 MAP-ELITES-ES

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

Table 5. MAP-ELITES hyperparameters

Parameter	Value
Number of centroids	1024
Total batch size b	256
GA batch size b_{GA}	256
Policy networks	[128, 128, $ \mathcal{A} $]
GA variation param. 1 σ_1	0.005
GA variation param. 2 σ_2	0.05

Table 6. MAP-ELITES-ES hyperparameters

Parameter	Value
Number of centroids	1024
Total batch size b	256
GA batch size b_{GA}	128
PG batch size b_{PG}	127
AI batch size b_{AI}	1
Policy networks	[128, 128, $ \mathcal{A} $]
GA variation param. 1 σ_1	0.005
GA variation param. 2 σ_2	0.05
Number of samples	1000
Sample sigma	0.02

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009