Solving Deep Reinforcement Learning Benchmarks with Linear Policy Networks

Abstract

Although Deep Reinforcement Learning (DRL) 1 methods can learn effective policies for challenging 2 problems such as Atari games and robotics tasks, 3 algorithms are complex and training times are often 4 long. This study investigates how evolution strate-5 6 gies (ES) perform compared to gradient-based deep 7 reinforcement learning methods. We use ES to optimize the weights of a neural network via neu-8 roevolution, performing direct policy search. We 9 benchmark both regular networks and policy net-10 works consisting of a single linear layer from ob-11 servations to actions; for three classical ES methods 12 and for three gradient-based methods such as PPO. 13 Our results reveal that ES can find effective lin-14 ear policies for many RL benchmark tasks, in con-15 trast to DRL methods that can only find success-16 ful policies using much larger networks suggest-17 ing that current benchmarks are easier to solve than 18 19 previously assumed. Interestingly, also for higher 20 complexity tasks, ES achieves results comparable to gradient-based DRL algorithms. Furthermore, 21 we find that by directly accessing the memory state 22 of the game, ES are able to find successful poli-23 cies in Atari, outperforming DQN. While gradient-24 based methods have dominated the field in recent 25 years, ES offers an alternative that is easy to imple-26 ment, parallelize, understand, and tune. 27

28 **1** Introduction

Gradient-based deep reinforcement learning (DRL) has 29 achieved remarkable success in various domains by en-30 abling agents to learn complex behaviors in challenging en-31 vironments based on their reward feedback, such as Star-32 Craft [Vinyals et al., 2019] and Go [Silver et al., 2016]. 33 However, new methods are often benchmarked on simpler 34 control tasks from OpenAI Gym, including the locomotion 35 tasks from MuJoCo [Haarnoja et al., 2018], or Atari games 36 [Mnih et al., 2015]. While it simplifies the comparison be-37 tween different approaches, these benchmarks may lack suf-38 ficient complexity, and performance may not always transfer 39 to more complicated tasks. Additionally, several studies have 40 indicated that DRL results are often hard to reproduce [Islam 41

et al., 2017], attributing these difficulties to the impact of the random seeds [Henderson *et al.*, 2018] and the choice of hyperparameters [Eimer *et al.*, 2023].

Evolution Strategies (ES) [Rechenberg, 1965; Bäck et al., 45 1991], a family of black-box optimization algorithms from 46 the field of Evolutionary Algorithms (EAs) [Bäck et al., 47 2023] have been studied as an alternative way to optimize 48 neural network weights, as opposed to conventional gradient-49 based backpropagation [Salimans et al., 2017; Such et al., 50 2017]. An ES is used to learn a controller for an RL task 51 by directly optimizing the neural network's weights, which 52 parameterize the RL policy. In this context, the ES is in-53 trinsically an RL method that performs direct policy search, 54 through neuroevolution [Igel, 2003]. In supervised learning, 55 gradient-based methods are often much more efficient than 56 ES for training NN weights, though more likely to be trapped 57 in local optima [Mandischer, 2002]. For RL, the need to 58 balance exploration with exploitation in gradient-based ap-59 proaches incurs more training steps to learn an optimal pol-60 icy [Igel, 2003], making ES an interesting alternative. While 61 EAs are not necessarily more sample-efficient, ES can be 62 more easily parallelized and scaled, offering the possibility 63 for faster convergence in terms of wall-clock time, and, being 64 a global search method, are less likely to get stuck in a local 65 optimum [Morse and Stanley, 2016]. 66

We benchmark three ES and three gradient-based RL meth-67 ods on well-known RL tasks to understand the circumstances 68 that favor ES over gradient-based methods. In particular, we 69 study the optimization of policy networks that consist of a 70 single linear layer, as they reduce the dimensionality of the 71 problem [Chrabaszcz et al., 2018; Rajeswaran et al., 2017]. 72 We compare these results to the larger networks that are used 73 by common gradient-based methods. Our main contributions 74 are as follows: 75

- ES can find effective linear policies for many RL benchmark tasks. In contrast, methods based on gradient descent need vastly larger networks.
- Contrary to the prevailing view that ES are limited to simpler tasks, they can address more complex challenges in MuJoCo. Gradient-based DRL only performs superiorly in the most challenging MuJoCo environments with more complex network architectures. This suggests that common RL benchmarks may be too simple or that con-



Figure 1: This study investigates how evolution strategies compare to gradient-based reinforcement learning methods in optimizing the weights of linear policies. We use both linear networks as the original DRL architectures to learn policies. We find that ES can learn linear policies for numerous tasks where DRL cannot, and in many instances, even surpasses the performance of the original DRL networks, such as in Swimmer.

86	plicated.
87	• Complex gradient-based approaches have dominated
88	DRL. However, ES can be equally effective, are algo-
89	rithmically simpler, allow smaller network architectures,
90	and are thus easier to implement, understand, and tune
91	(See Figure 1).

ventional gradient-based solutions may be overly com-

• We find that advanced self-adaptation techniques in ES are often not required for (single-layer) neuroevolution.

The rest of the paper is organized as follows: Section 2 discusses the background and related work of ES and DRL, our
algorithms are discussed in Section 3, the results are in Section 4, conclusions are in Section 5.

2 Background and Related Work

85

In RL, an agent learns from feedback through rewards and 99 penalties from its environment [Sutton and Barto, 2018]. 100 RL problems are formulated as a Markov Decision Process 101 $\langle S, A, P, R, \gamma \rangle$, where S is the set of states in the environ-102 ment, A the set of actions available to the agent, P the prob-103 ability of subsequent state transitions, R the reward function, 104 and $\gamma \in [0, 1]$ the discount factor [Bellman, 1957]. A policy 105 π can be computed of which action to take in each state. The 106 policy in DRL is typically represented by a deep neural net-107 work to map states to actions (probabilities). RL aims to find 108 the optimal policy π^* that maximizes the expected cumula-109 tive reward of a state. RL algorithms approach this goal in 110 different ways [Plaat, 2022]. The most common techniques 111 include value-function estimation [Watkins and Dayan, 1992; 112 Mnih et al., 2015], policy gradient methods [Williams, 1992], 113 actor-critic methods [Konda and Tsitsiklis, 1999; Schulman 114 et al., 2017; Haarnoja et al., 2018], and learning a model of 115 the environment [Hafner et al., 2020; Plaat et al., 2023]. 116

ES are a distinct class of evolutionary algorithms that are 117 particularly suitable for optimization problems in continuous 118 domains. ES begin with a population of randomly initial-119 ized candidate solutions in \mathbb{R}^n , with solutions represented 120 as *n*-dimensional vectors denoted by \mathbf{x} (like the policy) and 121 a given objective function $f : \mathbb{R}^n \to \mathbb{R}$ (like the reward). 122 Via perturbations using a parameterized multivariate normal 123 distribution, selection, and sometimes recombination, solu-124 tions evolve towards better regions in the search space [Bäck 125 et al., 1991]. Evolving neural networks with EAs is called 126 neuroevolution and can include the optimization of the net-127 work's weights, topology, and hyperparameters [Stanley et 128 al., 2019]. Using ES to evolve a neural network's weights is 129 similar to policy gradient methods in RL, where optimization 130 applies to the policy's parameter space. 131

Gradient-based deep RL has successfully tackled high-132 dimensional problems, such as playing video games 133 with deep neural networks encompassing millions parame-134 ters [Mnih et al., 2015; Vinyals et al., 2019]. However, 135 state-of-the-art ES variants are limited to smaller numbers of 136 parameters due to the computational complexity of, for ex-137 ample, the adaptation of the covariance matrix of the search 138 distribution. Covariance Matrix Adaptation Evolution Strat-139 egy (CMA-ES) is often used for dimensionality lower than 140 $n \leq 100$ [Müller and Glasmachers, 2018], and problems 141 with a dimensionality $n \ge 10.000$ become nearly impos-142 sible due to the memory requirements [Loshchilov, 2014]. 143 However, recent advancements have restricted the covari-144 ance matrix, in its simplest form, to its diagonal, to re-145 duce the computational complexity [Ros and Hansen, 2008; 146 Loshchilov, 2014; Nomura and Ono, 2022]. Others sample 147 from lower-dimensional subspaces [Maheswaranathan et al., 148 2019; Choromanski et al., 2019]. 149

In 2015, DRL reached a milestone by achieving superhuman performance in Atari games using raw pixel input [Mnih 151

et al., 2015]. This breakthrough marked a shift in RL to-152 wards more complicated, high-dimensional problems and a 153 shift from tabular to deep, gradient-based methods. For sim-154 pler tasks, the CMA-ES has been used to evolve neural net-155 works for pole-balancing tasks, benefiting from covariance 156 matrix to find parameter dependencies, enabling faster op-157 timization [Igel, 2003; Heidrich-Meisner and Igel, 2009]. 158 While the use of evolutionary methods for RL can be traced 159 back to the early 90s [Whitley et al., 1993; Moriarty and 160 Mikkulainen, 1996], the paper by [Salimans et al., 2017] 161 rekindled interest in ES from the field of RL as an alterna-162 tive for gradient-based methods in more complicated tasks. 163 Researchers showed that a natural evolution strategy (NES) 164 can compete with deep RL in robot control in MuJoCo and 165 Atari games due to its ability to scale over parallel workers. 166 In contrast to deep methods where entire gradients are typi-167 cally shared, the workers only communicate the fitness score 168 and the random seed to generate the stochastic perturbations 169 of the policy parameters. Studies have subsequently demon-170 strated that simpler methodologies can yield results compara-171 ble to NES, such as a classical ES [Chrabaszcz et al., 2018] 172 and Augmented Random Search [Mania et al., 2018], which 173 closely resembles a global search heuristic from the 1990s 174 [Salomon, 1998]. In addition, when separating the computer 175 vision task from the actual policy in playing Atari, the size 176 of the neural network can be drastically decreased [Cuccu 177 et al., 2018], and policies with a single linear layer, map-178 ping states directly to actions, can effectively solve the con-179 tinuous control tasks [Chrabaszcz et al., 2018; Rajeswaran 180 et al., 2017]. The development of dimension-lowering tech-181 niques, such as world models [Ha and Schmidhuber, 2018; 182 Hafner et al., 2020] and autoencoders [Hinton and Salakhut-183 dinov, 2006], also opens up new possibilities for ES to effec-184 tively solve more complex problems by simplifying them into 185 more manageable forms. 186

3 Methods 187

We implement three ES methods and benchmark them against 188 three popular gradient-based DRL methods. 189

3.1 Gradient-Based Algorithms 190

We use three popular gradient-based DRL algorithms: Deep 191 O-learning [Mnih et al., 2015], Proximal Policy Optimiza-192 tion [Schulman et al., 2017], and Soft Actor-Critic [Haarnoja 193 et al., 2018]. We summarize the main gradient-update ideas 194 below. 195

Deep Q-Learning 196

Deep Q-learning (DQN) combines a deep neural network 197 with Q-learning to learn the value function in a high-198 dimensional environment [Mnih et al., 2015]. Each experi-199 ence tuple (s_t, a_t, r_t, s_{t+1}) is stored in a replay buffer. The 200 agent randomly selects a batch of experiences to update the 201 value function. The replay buffer breaks the correlation be-202 tween consecutive experiences, leading to lower variance. 203 The primary Q-network weights θ are updated every training 204 205 step by minimizing the expectation of the squared difference between the predicted Q-value of the current state-action pair 206

$$Q(s, a; \theta)$$
 and the target Q-value $Q(s', a'; \theta^{-})$:

$$L(\theta) = \mathbb{E}\left[\left(r + \gamma \max_{a'} Q(s', a'; \theta^{-}) - Q(s, a; \theta)\right)^{2}\right]$$

The weights from the primary Q-network are copied every 208 N timesteps to a separate target network $\theta^- \leftarrow \theta$ to prevent 209 large oscillations in the loss function's target values. 210

Proximal Policy Optimization

Proximal Policy Optimization (PPO) was introduced to 212 improve the complexity of earlier policy gradient meth-213 ods [Schulman et al., 2017]. PPO introduces a simpler, 214 clipped objective function: 215

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta) \hat{A}_t, \operatorname{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

where $\hat{\mathbb{E}}_t$ denotes the empirical expectation over a finite batch 216 of samples, the probability ratio $r_t(\theta)$ reflects the probability 217 of an action under the current policy compared to the previous 218 policy, \hat{A}_t is the advantage estimate at timestep t, and ϵ is 219 a hyperparameter defining the clipping range. The clipping 220 mechanism clips the ratio $r_t(\theta)$ within the range $[1-\epsilon, 1+\epsilon]$. 221

Soft Actor-Critic

SAC objective's function maximizes the expected return and 223 entropy simultaneously to ensure a balanced trade-off be-224 tween exploitation and exploration: 225

$$\pi^* = \arg \max_{\pi} \sum_{t} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} \left[r(s_t, a_t) + \alpha H(\pi(\cdot | s_t)) \right] ,$$

where α is the temperature parameter that scales the importance of the entropy $H(\pi(\cdot|s_t))$ of the policy π given the state s_t . SAC updates its Q-value estimates using a soft Bellman backup operation that includes an entropy term:

$$\begin{aligned} Q_{\text{new}}(s_t, a_t) &= \mathbb{E}_{s_{t+1} \sim \mathcal{E}} \left[r(s_t, a_t) \right. \\ &+ \gamma \left(Q_{\text{old}}(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1} | s_{t+1}) \right) \right]. \end{aligned}$$

SAC employs twin Q-networks to mitigate overestimation 226 bias and stabilize policy updates by using the minimum of 227 their Q-value estimates. 228

3.2 Evolution Strategies

ES are designed for solving continuous optimization prob-230 lems maximize_x $f(\mathbf{x})$, where $f : \mathbb{R}^n \to \mathbb{R}$. The ES meth-231 ods are used here to optimize the neural network parameters 232 for the policy function, i.e., the ES is used for neuroevolu-233 tion of the weights of the neural network that represents the 234 policy. The methods are variants of derandomized ES and 235 use a parameterized normal distribution $\mathcal{N}(\mathbf{m}^{(g)}, \sigma^{(g)}\mathbf{C}^{(g)})$ 236 to control the direction of the search. The algorithm adapts 237 the parameters of the search distribution to achieve fast con-238 vergence (Algorithm 1). 239

The ES samples λ offspring from its mutation distribution at each iteration. By selecting the $\mu < \lambda$ most promising offspring to update its parameters, it moves to regions of higher optimality. After sorting the μ offspring by fitness ranking, the mean of the search distribution is updated via weighted recombination:

$$\mathbf{m}^{(g+1)} = \mathbf{m}^{(g)} + c_m \sum_{i=1}^{\mu} w_i(\mathbf{x}_i - \mathbf{m}^{(g)})$$

211

222

229

Algorithm 1 Generic ES

Require: Objective function f, number of offspring λ , number of parents μ , initial estimates for $\mathbf{m}^{(0)}$ and $\sigma^{(0)}$ $\mathbf{C}^{(0)} \leftarrow \mathbf{I}_n$ **for** g in 1, 2, ... **do** Sample λ candidates $\mathbf{x}_i \sim \mathcal{N}(\mathbf{m}^{(g)}, \sigma^{(g)}\mathbf{C}^{(g)})$ Evaluate objective function $f_i \leftarrow f(\mathbf{x}_i)$ Select and rank μ best candidates Adapt $\mathbf{m}^{(g+1)}, \sigma^{(g+1)}, \mathbf{C}^{(g+1)}$ **end for**

The ES variants adapt, with increasing complexity, the 240 scale $\sigma^{(g)}$ and shape $\mathbf{C}^{(g)}$ of the mutation distribution. The 241 Cumulative Step-size Adaptation Evolution Strategy (CSA-242 ES) only adapts $\sigma^{(g)}$, producing exclusively isotropic (i.e. 243 $\mathbf{C}^{(g)} = \mathbf{I}_n$) mutations during optimization. The separa-244 ble Covariance Matrix Adaptation Evolution Strategy (sep-245 CMA-ES) additionally adapts the diagonal entries of the co-246 variance matrix $\mathbf{C}^{(g)}$, producing mutation steps with arbitrary 247 scaling parallel to each coordinate axis. Finally, the CMA-ES 248 adapts the full covariance matrix, which allows the mutation 249 distribution to scale to arbitrary rotations. Figure 2 illustrates 250 the evolution of the mutation distribution for each of these 251 three methods when optimizing a two-dimensional quadratic 252 function. The figure shows that the mutation distribution 253 guides the search, favoring selected mutation steps with high 254 probability [Hansen and Ostermeier, 2001]. The CSA-ES 255 uses a process called cumulation of historically selected mu-256 tation steps in order to update the value of the global step 257 size parameter $\sigma^{(g)}$. We implemented the algorithm follow-258 ing [Chotard et al., 2012], using recommended hyperparame-259 ter settings. While several modifications of the CMA-ES have 260 been developed over the years, we implemented a canonical 261 version of the algorithm, as first introduced in [Hansen and 262 Ostermeier, 2001]. The update of the full covariance ma-263 trix becomes computationally impractical for n > 100, but 264 the sep-CMA-ES, which we implemented according to [Ros 265 and Hansen, 2008], does not suffer from this restriction. As 266 267 shown in Figure 2, this algorithm only computes variances for each coordinate axis, which makes it applicable to much 268 higher dimensions, as the computational complexity for the 269 update of the mutation distribution scales only linearly with 270 271 n.

272 **3.3** Network Architecture

We compare the performance of linear policies trained 273 through neuroevolution by ES with gradient-based methods 274 inspired by earlier studies demonstrating this approach's fea-275 sibility [Mania et al., 2018; Rajeswaran et al., 2017]. For the 276 ES, only linear policies are trained, defined as a linear map-277 ping from states to actions, activated by either an argmax 278 or tanh function for discrete and continuous action spaces, 279 respectively (no hidden layer: a fully connected shallow net-280 work). We use the gradient-based methods to train the same 281 linear policies for each control task. Table 2 (Supplementary 282 283 material) shows the number of trainable weights for each environment for a linear policy. Additionally, a network archi-284



Figure 2: Adaptation of the mutation distribution for three different Evolution Strategies for the first ten generations of a twodimensional quadratic function. Function values are shown with color; darker indicates lower (better). Top row: mutation distribution for CSA-ES; middle row: sep-CMA-ES; bottom row: CMA-ES

tecture based on the original studies for each of the gradient-285 based methods is trained for comparison. We employ the ar-286 chitecture from the original studies for PPO [Schulman et al., 287 2017] and SAC [Haarnoja et al., 2018]. For DQN, we use the 288 default architecture from CleanRL's library, which has been 289 tested across multiple control environments and showed good 290 results [Huang et al., 2022]. Specifics for these architectures 291 and other hyperparameters can be found in the supplementary 292 material. We do not train these deep architectures using ES, 293 as they only serve as a benchmark to demonstrate the intended 294 usage of the gradient-based algorithms, and self-adaptation 295 mechanisms are increasingly less useful for such high dimen-296 sions [Chrabąszcz et al., 2018]. 297

298

3.4 Experimental Setup

We conduct experiments on common control tasks of varying 299 complexity levels from the Gymnasium API [Towers et al., 300 2023]. For each of the considered environments, five runs 301 using different random seeds are conducted for each algo-302 rithm/control task combination to test the stability of each ap-303 proach. Value-based DQN is only used for environments with 304 discrete action spaces, SAC for continuous action spaces, and 305 PPO and ES are used for both action spaces. The RL algo-306 rithms are implemented using the cleanRL library¹ that has 307 been benchmarked across several environments; we removed 308 the hidden layers for the linear network. The specifics of 309 the Evolution Strategies (ES) implementations are detailed in 310 the repository, the link to which is provided in the footnote². 311 Since the environments are stochastic, we report the median 312 episodic return, calculated over five test episodes. As was dis-313 cussed in [Salimans et al., 2017], for ES, the wall-clock time 314 required to solve a given control task decreases linearly with 315 the number of parallel workers available. This allows us to 316

¹https://github.com/vwxyzjn/cleanrl/tree/master

²anonymized for reviewing



Figure 3: Training curves for the CartPole, LunarLander, Swimmer, HalfCheetah, Boxing, and SpaceInvaders environments. Episodic return (calculated using 5 test episodes) vs. the number of training timesteps is shown. Each curve represents the median of 5 trial runs conducted with different random seeds; the shaded area denotes standard deviations. The results show that ES solve the classic control environments. Cartpole and LunarLander almost immediately, surpassing the gradient descent methods. Even for the more difficult Swimmer environment, ES find a linear policy outperforming DRL in terms of timesteps and performance. While SAC outperforms all other methods in Cheetah, linear ES outperforms classic PPO. For the Atari environments, Boxing and Space Invaders, ES is able to learn a linear policy from the RAM input, while linear DQN fails to do so, and only for Boxing deep DQN finds a successful policy.

Algorithm	CSA-ES	CMA-ES	sep-CMA	Random	Human	DQN	ES*
Atlantis	84690	87100	88580	12850	29028	85641	103500
B. Rider	2215	1967	2222	363.9	5775	6846	5072
Boxing	96.0	96.8	95.1	0.1	4.3	71.8	100
C. Climber	36170	29290	32940	10781	35411	114103	57600
Enduro	65.1	58.9	69.0	0	309.6	301.8	102
Pong	5.7	7.4	7.1	-20.7	9.3	18.9	21
Q*Bert	7355	5732	7385	163.9	13455	10596	14700
Seaquest	959	948	954	68.4	20182	5286	1470
S. Invaders	1640	1972	1488	148	1652	1976	2635

Table 1: Average maximum score per game across trials. The rightmost column shows the best-performing ES episode per game. For comparison, the scores for DQN, a random agent, and a Human player taken from [Mnih *et al.*, 2015] are shown. The highest scores are shown in boldface. The results show that an ES outperforms the original DQN scores for Atlantis and Boxing. For the other games, the highest score is attained by the DQN agent, although CMA-ES achieves a score almost identical to DQN on SpaceInvaders. Furthermore, the best ES policy often matches the performance of DQN, demonstrating that a linear policy can be just as effective. Standard deviation values are provided in Table 1 in the supplementary material.

do many more evaluations of the environment than is feasible with gradient-based RL. For fairness of comparison, we limit the difference in the number of training time steps allowed by a single order of magnitude. Specific hyper-parameters used for each environment can be found in the supplementary material, including hardware. For the ES, we initialize each experiment with $\mathbf{m}^{(0)} = \mathbf{0}$. We calculate a rolling mean and variance of the observations of the environment during each run. These values are then used to normalize each state observation to standard normal entries by subtracting this rolling mean and dividing by the standard deviation.

328 Classic RL Environments

The first set of experiments includes the classic control tasks Cartpole, Acrobot, and Pendulum. We include Bipedal-Walker and LunarLander from the Box2D simulations for slightly more complex dynamics. Each run uses a maximum of 500 000 timesteps for each environment. The exception to this is the BipedalWalker task, for which $2 \cdot 10^6$ timesteps are used.

336 MuJoCo Simulated Robotics

We evaluate the algorithms on the MuJoCo environments 337 [Todorov *et al.*, 2012] for higher complexity levels, includ-338 ing Hopper, Walker2D, HalfCheetah, Ant, Swimmer, and Hu-339 manoid. Table 3 (Supplementary material) provides training 340 details. As was noted by [Mania et al., 2018], ES methods 341 have exploration at the policy level, whereas gradient-based 342 methods explore on the action level. In the locomotion tasks, 343 a positive reward is provided for each time step where the 344 agent does not fall over. This causes the ES method to stay 345 in a local optimum when the agent stavs upright but does not 346 move forward (the gradient-based methods do not get stuck). 347 Following [Mania et al., 2018], we modified the reward func-348 tion for these environments for ES, subtracting the positive 349 stability bonus and only rewarding forward locomotion. 350

351 Atari Learning Environment

Finally, we benchmark DQN against the ES with linear poli-352 cies on games from the Atari suite. To demonstrate the effec-353 tiveness of linear policies for these high-dimensional tasks, 354 we take inspiration from the approach by [Cuccu et al., 2018] 355 and separate the computer vision task from the control task. 356 We train the ES agents on the 128 bytes of the Random Ac-357 cess Memory (RAM) in the simulated Atari console. This 358 drastically reduces the input dimensionality of the controller, 359 allowing for the training of smaller policies. This assumes 360 that the random access memory sufficiently encodes the state 361 of each game without having to extract the state from the raw 362 pixel images. It should be noted that not for all games is 363 RAM information sufficient to train a controller and that for 364 365 some games, DQN is easier to train on pixel images than on RAM input [Sygnowski and Michalewski, 2017]. Since we 366 evaluate linear policies, we fix frame skipping to 4, with no 367 sticky actions [Machado et al., 2018], similar to the settings 368 used in [Mnih et al., 2015]. For each run, the ES were trained 369 for 20000 episodes, where the maximum episode length was 370 capped at 270 000 timesteps. The gradient-based methods 371 were trained for a maximum of $2 \cdot 10^7$ timesteps. 372

373 **4 Results**

In this section, we present the results of our experiments (more details, including training curves and tables with summary statistics, can be found in the Supplementary material). Here, we focus on a selection of environments (Figure 3).

378 4.1 Classics RL Environments

The first column of Figure 3 shows the training curves for the gradient-based methods and for ES on the CartPole and LunarLander environments (the Supplementary material has more). For both environments, the ES policies outperform the deep gradient-based methods with just a simple linear 383 policy. For CartPole, the results are even more surprising, 384 as the ES policies are able to solve the environment in the 385 first few iterations of training through pure random sampling 386 from a standard normal distribution. The deep gradient-based 387 methods, on the other hand, require around $2 \cdot 10^5$ timesteps 388 in order to solve CartPole. The gradient-based methods are 389 unable to find a good linear policy for CartPole. This pat-390 tern persists for LunarLander, where, even though the ES re-391 quires around $2 \cdot 10^5$ timesteps to solve the environment, the 392 gradient-based methods cannot find a good linear policy at all. 393 While the deep gradient-based methods eventually seem to 394 catch up to the ES, it still requires more than $5 \cdot 10^5$ timesteps 395 to train a stable policy for LunarLander. For the Bipedal-396 Walker task (see Supplementary material), DQN and PPO are 397 able to find good policies, with the gradient-based methods 398 requiring fewer timesteps. Only SAC is able to solve Pen-399 dulum within $5 \cdot 10^5$ timesteps. For Acrobot, again, the ES 400 are able to solve the environment almost instantly, while the 401 gradient-based methods require a good number of environ-402 ment interactions to do so. 403

404

427

4.2 MuJoCo Simulated Robotics

The center column of Figure 3 shows that the ES policies 405 are much better at finding a policy for the Swimmer environ-406 ment, while for HalfCheetah, SAC greatly outperforms all 407 other methods. However, ES outperforms PPO. Moreover, 408 none of the gradient-based methods are able to find a good 409 linear policy for HalfCheetah and Swimmer. This pattern 410 holds for almost all MuJoCo experiments; only for the Ant 411 environment can PPO find a linear policy with decent per-412 formance. Overall, as the number of weights increases, the 413 performance of the ES lags behind the deep gradient-based 414 methods (see Table 2 in the Supplementary material). Nev-415 ertheless, even though ES generally require more timesteps, 416 they can still find good linear policies for most environments, 417 which are just as effective as policies found by vastly larger 418 networks (see Table 4 in the Supplementary material). Even 419 for the most complex of these environments, Humanoid, the 420 ES are able, in several trials, to find a linear policy that has a 421 higher episodic return, ≈ 8000 (averaged over 5 test episodes 422 with different random seeds), than was found by the best deep 423 gradient-based method, SAC. Furthermore, ES timesteps are 424 quicker and easier to parallelize, meaning experiments are 425 much faster to run. 426

4.3 Atari Learning Environment

The last column of Figure 3 shows the training curves of 428 ES vs. DQN for the Atari games SpaceInvaders and Box-429 ing. The figure shows that when training agents that use the 430 controller's RAM state as observations, ES outperform linear 431 DQN in most cases. CrazyClimber (Supplementary material) 432 is the only exception, for which the performance of ES and 433 linear DQN is similar. Even comparing against deep DQN 434 trained on RAM memory, we find that for both the games in 435 Figure 3, the ES yield better policies and require fewer envi-436 ronment interactions. In addition, Table 1 shows the average 437 highest score per trial for each of the tested games for the ES, 438 compared against the numerical results presented in [Mnih et 439

al., 2015] for a Human, Random and a DQN player that uses 440 pixel input. The table additionally shows the highest score 441 attained by any ES in any trial, averaged over 5 test episodes. 442 For both Atlantis and Boxing, an ES achieves the highest 443 score. For all the other games tested, the DQN agent earned 444 a higher score than all RAM-trained ES, although CMA-ES 445 achieved a score almost identical to DQN on SpaceInvaders. 446 This score is attained by an agent that uses a linear policy con-447 sisting of only 768 weights, while the policy trained by DQN 448 has $\approx 1.5 \cdot 10^6$ (pixel-based, deep policy network). Moreover, 449 the best-found policy by any ES is often competitive with 450 DQN. This indicates that a linear policy, which is competi-451 tive with pixel-based DQN, does exist. 452

453 **5** Discussion and Conclusion

In this study, we have studied different ways to optimize rein-454 forcement learning policies with conventional deep learning 455 gradient-based backpropagation methods as used in DQN, 456 PPO, and SAC, and with three evolution strategy methods 457 (ES). We have applied these methods to several classic re-458 inforcement learning benchmarks. For these methods, we 459 trained the regular deep network as conventionally used, as 460 well as a neural network with no hidden layers, i.e., a lin-461 ear mapping from states to actions, as a low-complexity con-462 troller for each environment. In many of the tested envi-463 ronments, the linear policies trained with the ES are on par 464 465 or, in some cases, even better controllers than the deep policy networks trained with the gradient-based methods. Ad-466 ditionally, the gradient-based methods are often ineffective 467 at training these simple policies, requiring much deeper net-468 works. For our experiments on Atari, we find that by access-469 ing the RAM memory of the Atari controller, ES methods 470 can find a linear policy that is competitive with "superhuman" 471 DQN [Mnih et al., 2015]. It should be noted that there are cer-472 tain high-complexity environments where the deep gradient-473 based methods yield better policies, e.g. SAC for HalfChee-474 tah. However, even for these environments, linear policies 475 exist that are competitive and much more easily interpretable. 476 We conclude that conventional gradient-based methods 477 might be overly complicated or that more complex bench-478 marks are required to properly evaluate algorithms. In fact, 479 even for the most complex locomotion task included in our 480 experiments, the Humanoid environment, the CMA-ES was 481 able to find a linear policy that was competitive with state-482 of-the-art methods. As the ES are stochastic algorithms, they 483 were not able to find these policies for every trial run, but our 484 results show that such policies do exist. We expect the search 485 landscapes for these environments to be deceptive and mul-486 timodal, and future work could help discover effective algo-487 rithms for more consistently training these linear policy net-488 works, for example, using niching methods [Shir and Bäck, 489 2005]. We hypothesize that gradient-based methods may 490 struggle to find linear policies due to the multimodal nature 491 of the search landscape, a phenomenon also seen in super-492 vised learning [Kawaguchi, 2016]. Counterintuitively, with 493 gradient-based methods, it seems more straightforward to 494 495 train deeper architectures than shallower ones with far fewer weights, as was also shown by [Schwarzer et al., 2023]. We 496

note that gradient-based methods are essentially local search 497 methods, requiring heuristically controlled exploration, while 498 ES, in the early phases of optimization, are performing global 499 search, producing more diverse solutions. This also becomes 500 evident in the more simple environments, such as CartPole, 501 where the ES are able to almost instantly sample the optimal 502 policy, while the gradient-based methods have a much harder 503 time. 504

Moreover, we find many counterexamples to the prevailing 505 view that ES are less sample efficient than the gradient-based 506 methods. For many of the low to medium-complexity envi-507 ronments, the ES are actually more sample efficient and re-508 quire fewer environment interactions than the deep gradient-509 based methods. On the other hand, for the more complex 510 environments, and with increasing dimensionality, we find 511 that the ES can take more time steps to converge than the 512 deep gradient-based methods. This is to be expected, as the 513 self-adaptation mechanisms that are central to the ES become 514 increasingly ineffective for larger dimensions [Chrabaszcz et 515 al., 2018; Müller and Glasmachers, 2018]. We have com-516 pared three ES, that, with increasing levels of complexity, 517 adapt the shape of the mutation distribution in order to con-518 verge the search. Our results indicate that updates of the 519 covariance matrix are often not required and that perform-520 ing step size adaptation is sufficient. While we expect the 521 search landscape to be multimodal, relative scaling and ro-522 tation of search coordinates seem absent, allowing isotropic 523 mutations to be effective for these problems. This would 524 also explain the effectiveness of the approach demonstrated 525 in [Maheswaranathan et al., 2019], which would be heav-526 ily impacted by conditioning on the search space. However, 527 this may be explained because optimizing a single linear layer 528 may exhibit less inherent variance and covariance than multi-529 ple layers. 530

Overall, we have demonstrated the potential of linear poli-531 cies on popular RL benchmarks. We showed that ES are 532 effective optimizers for these policies compared to gradient-533 based methods. Additionally, we note that ES are simpler in 534 design, have fewer hyperparameters, and are trivially paral-535 lelizable. Hence, ES can perform more environment inter-536 actions within the same time frame [Salimans et al., 2017]. 537 Moreover, evaluating linear policies is faster than evaluat-538 ing one or sometimes several deep architectures, making the 539 training much more expedient regarding wall-clock time. As 540 the need for energy-efficient policy networks increases, our 541 results warrant a closer look at ES for RL and training of sim-542 pler policies for tasks currently considered complex. For fu-543 ture work, we aim to extend our benchmark with more types 544 of classical ES and strategies for multimodal optimization 545 [Preuss, 2015]. Additionally, it would be interesting to study 546 the effect of the step-size adaptation methods in the presence 547 of one or more hidden layers. Next, we will explore the po-548 tential of linear networks for other applications. Inspired by 549 works such as [Cuccu et al., 2018], we will look at more com-550 plex Atari games to see if they can also be solved by simple, 551 energy-efficient means. 552

553 **References**

- [Bäck *et al.*, 1991] T Bäck, F Hoffmeister, and H-P Schwe fel. A survey of evolution strategies. In *Proceedings of the fourth international conference on genetic algorithms*.
 Citeseer, 1991.
- [Bäck *et al.*, 2023] T Bäck, A V Kononova, B van Stein,
 H Wang, K Antonov, R Kalkreuth, J de Nobel, D Vermetten, R de Winter, and F Ye. Evolutionary algorithms for
- parameter optimization—thirty years later. *Evolutionary Computation*, 31(2):81–122, 2023.
- [Bellman, 1957] R Bellman. A markovian decision process.
 Journal of Mathematics and Mechanics, pages 679–684, 1957.
- [Choromanski *et al.*, 2019] K M Choromanski, A Pacchiano, J Parker-Holder, Y Tang, and V Sindhwani. From complexity to simplicity: Adaptive es-active subspaces for blackbox optimization. *Advances in Neural Information*
- 570 *Processing Systems*, 32, 2019.
- ⁵⁷¹ [Chotard *et al.*, 2012] A Chotard, A Auger, and N Hansen.
 ⁵⁷² Cumulative step-size adaptation on linear functions. In
- 573 Parallel Problem Solving from Nature-PPSN XII: 12th
- 574 International Conference, Taormina, Italy, September 1-
- 575 *5, 2012, Proceedings, Part I 12*, pages 72–81. Springer, 2012.
- 577 [Chrabąszcz *et al.*, 2018] P Chrabąszcz, I Loshchilov, and
 F Hutter. Back to basics: Benchmarking canonical evo579 lution strategies for playing atari. In *Proceedings of the*580 *Twenty-Seventh International Joint Conference on Arti-*581 *ficial Intelligence, IJCAI-18*, pages 1419–1426. Interna-
- tional Joint Conferences on Artificial Intelligence Organization, 2018.
- [Cuccu *et al.*, 2018] G Cuccu, J Togelius, and P Cudré Mauroux. Playing atari with six neurons. *arXiv preprint arXiv:1806.01363*, 2018.
- [Eimer *et al.*, 2023] T Eimer, M Lindauer, and R Raileanu.
 Hyperparameters in reinforcement learning and how to
 tune them. In *International conference on machine learn- ing*, 2023.
- ⁵⁹¹ [Ha and Schmidhuber, 2018] D Ha and J Schmidhuber.
 ⁵⁹² World models. *arXiv preprint arXiv:1803.10122*, 2018.
- ⁵⁹³ [Haarnoja *et al.*, 2018] T Haarnoja, A Zhou, K Hartikainen,
- G Tucker, S Ha, J Tan, V Kumar, H Zhu, A Gupta,
 P Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- [Hafner *et al.*, 2020] D Hafner, T Lillicrap, J Ba, and
 M Norouzi. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.
- [Hansen and Ostermeier, 2001] N Hansen and A Ostermeier.
 Completely derandomized self-adaptation in evolution
 strategies. *Evolutionary computation*, 9(2):159–195,
 2001.
- ⁶⁰⁵ [Heidrich-Meisner and Igel, 2009] V Heidrich-Meisner and
 ⁶⁰⁶ C Igel. Neuroevolution strategies for episodic reinforce-

ment learning. Journal of Algorithms, 64(4):152–168, 607 2009.

- [Henderson *et al.*, 2018] P Henderson, R Islam, P Bachman,
 J Pineau, D Precup, and D Meger. Deep reinforcement
 learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [Hinton and Salakhutdinov, 2006] G Hinton and R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006. 615
- [Huang *et al.*, 2022] S Huang, R Fernand, J Dossa, C Ye,
 J Braga, P Chakraborty, K Mehta, and J Araújo. Cleanrl:
 High-quality single-file implementations of deep rein forcement learning algorithms. *Journal of Machine Learn- ing Research*, 23(274):1–18, 2022.
- [Igel, 2003] C Igel. Neuroevolution for reinforcement learning using evolution strategies. In *The 2003 Congress* on Evolutionary Computation, 2003. CEC'03., volume 4, pages 2588–2595. IEEE, 2003.
- [Islam et al., 2017] R Islam, P Henderson, M Gomrokchi, and D Precup. Reproducibility of benchmarked deep reinforcement learning tasks for continuous control. In *Reproducibility in Machine Learning Workshop (ICML)*, 2017.
- [Kawaguchi, 2016] K Kawaguchi. Deep learning without poor local minima. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [Konda and Tsitsiklis, 1999] V Konda and J Tsitsiklis.
 Actor-critic algorithms. Advances in neural information
 processing systems, 12, 1999.
 636
- [Loshchilov, 2014] I Loshchilov. A computationally efficient
 limited memory cma-es for large scale optimization. In
 Proceedings of the 2014 Annual Conference on Genetic
 and Evolutionary Computation, pages 397–404, 2014.
- [Machado *et al.*, 2018] M Machado, M Bellemare, E Talvitie, J Veness, M Hausknecht, and M Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018. 645
- [Maheswaranathan *et al.*, 2019] N Maheswaranathan,
 L Metz, G Tucker, D Choi, and J Sohl-Dickstein. Guided
 evolutionary strategies: Augmenting random search with
 surrogate gradients. In *International Conference on Machine Learning*, pages 4264–4273. PMLR, 2019.
- [Mandischer, 2002] M Mandischer. A comparison of evolution strategies and backpropagation for neural network training. *Neurocomputing*, 42(1):87–117, 2002. 653
- [Mania *et al.*, 2018] H Mania, A Guy, and B Recht. Simple random search provides a competitive approach to reinforcement learning. *arXiv preprint arXiv:1803.07055*, 656 2018. 657
- [Mnih et al., 2015] V Mnih, K Kavukcuoglu, D Silver, A A
 Rusu, J Veness, M G Bellemare, A Graves, M Ried miller, A K Fidjeland, G Ostrovski, et al. Human-level

- control through deep reinforcement learning. Nature. 661 518(7540):529-533, 2015. 662
- [Moriarty and Mikkulainen, 1996] D Moriarty and 663 R Mikkulainen. Efficient reinforcement learning through 664 symbiotic evolution. Machine learning, 22:11-32, 1996. 665
- [Morse and Stanley, 2016] G Morse and K Stanley. Simple 666
- evolutionary optimization can rival stochastic gradient de-667 scent in neural networks. In Proceedings of the Genetic 668
- and Evolutionary Computation Conference 2016, pages 669 477-484, 2016. 670
- [Müller and Glasmachers, 2018] N Müller and T Glasmach-671 ers. Challenges in high-dimensional reinforcement learn-672
- ing with evolution strategies. In Parallel Problem Solving 673
- from Nature–PPSN XV: 15th International Conference, 674
- Coimbra, Portugal, September 8–12, 2018, Proceedings, 675 Part II 15, pages 411-423. Springer, 2018. 676
- [Nomura and Ono, 2022] M Nomura and I Ono. Fast mov-677
- ing natural evolution strategy for high-dimensional prob-678
- lems. In 2022 IEEE Congress on Evolutionary Computa-679 tion (CEC), pages 1-8. IEEE, 2022. 680
- [Plaat et al., 2023] A Plaat, W Kosters, and M Preuss. High-681 accuracy model-based reinforcement learning, a survey. 682 Artificial Intelligence Review, pages 1–33, 2023. 683
- [Plaat, 2022] A Plaat. Deep Reinforcement Learning. 684 Springer Verlag, Singapore, 2022. 685
- [Preuss, 2015] M Preuss. Multimodal Optimization by 686 Means of Evolutionary Algorithms. Springer Publishing 687 Company, Incorporated, 1st edition, 2015. 688
- [Rajeswaran et al., 2017] A Rajeswaran, Κ Lowrey, 689 E Todorov, and S Kakade. Towards generalization and 690 simplicity in continuous control. Advances in Neural 691 Information Processing Systems, 30, 2017. 692
- [Rechenberg, 1965] I Rechenberg. Cybernetic solution path 693 of an experimental problem. Royal Aircraft Establishment 694 Library Translation 1122, 1122, 1965. 695
- [Ros and Hansen, 2008] R Ros and N Hansen. A simple 696 modification in cma-es achieving linear time and space 697 complexity. In International conference on parallel prob-698 lem solving from nature, pages 296–305. Springer, 2008. 699
- [Salimans et al., 2017] T Salimans, J Ho, X Chen, S Sidor, 700 and I Sutskever. Evolution strategies as a scalable 701 alternative to reinforcement learning. arXiv preprint 702 arXiv:1703.03864, 2017. 703
- [Salomon, 1998] R Salomon. Evolutionary algorithms and 704 gradient search: similarities and differences. IEEE Trans-705 actions on Evolutionary Computation, 2(2):45-55, 1998. 706
- [Schulman et al., 2017] J Schulman, F Wolski, P Dhariwal, 707
- A Radford, and O Klimov. Proximal policy optimization 708 algorithms. arXiv preprint arXiv:1707.06347, 2017. 709
- [Schwarzer et al., 2023] M Schwarzer, J Ceron, A Courville, 710 M Bellemare, R Agarwal, and P Castro. Bigger, bet-711
- ter, faster: Human-level atari with human-level efficiency. 712
- In International Conference on Machine Learning, pages 713 30365-30380. PMLR, 2023. 714

- [Shir and Bäck, 2005] O Shir and T Bäck. Niching in evo-715 lution strategies. In Proceedings of the 7th annual con-716 ference on Genetic and evolutionary computation, pages 717 915-916, 2005. 718
- [Silver et al., 2016] D Silver, A Huang, C J Maddison, 719 A Guez, L Sifre, G Van Den Driessche, J Schrittwieser, 720 I Antonoglou, V Panneershelvam, M Lanctot, et al. Mas-721 tering the game of go with deep neural networks and tree 722 search. Nature, 529(7587):484-489, 2016. 723
- [Stanley et al., 2019] K O Stanley, J Clune, J Lehman, and 724 R Miikkulainen. Designing neural networks through neu-725 roevolution. *Nature Machine Intelligence*, 1(1):24–35, 726 2019. 727
- [Such et al., 2017] F P Such, V Madhavan, E Conti, 728 J Lehman, K O Stanley, and J Clune. Deep neuroevolu-729 tion: Genetic algorithms are a competitive alternative for 730 training deep neural networks for reinforcement learning. 731 arXiv preprint arXiv:1712.06567, 2017. 732
- [Sutton and Barto, 2018] R S Sutton and A G Barto. Rein-733 forcement learning: An introduction. MIT press, 2018. 734
- [Sygnowski and Michalewski, 2017] J Sygnowski and 735 H Michalewski. Learning from the memory of atari 2600. 736 In Computer Games: 5th Workshop on Computer Games, 737 CGW 2016. and 5th Workshop on General Intelligence in 738 Game-Playing Agents, GIGA 2016, Held in Conjunction 739 with the 25th International Conference on Artificial 740 Intelligence, IJCAI 2016, New York, USA, July 9-10, 741 2016, Revised Selected Papers 5, pages 71-85. Springer, 742 2017. 743
- [Todorov et al., 2012] E Todorov, T Erez, and Y Tassa. Mu-744 joco: A physics engine for model-based control. In 2012 745 IEEE/RSJ international conference on intelligent robots 746 and systems, pages 5026-5033. IEEE, 2012. 747
- [Towers et al., 2023] M Towers, J Terry, A Kwiatkowski, 748 J Balis, G de Cola, T Deleu, M Goulão, A Kallinteris, 749 A KG, M Krimmel, R Perez-Vicente, A Pierré, S Schul-750 hoff, J J Tai, A T J Shen, and O Younis. Gymnasium, 751 March 2023. 752
- [Vinyals et al., 2019] O Vinyals, I Babuschkin, W M Czar-753 necki, M Mathieu, A Dudzik, J Chung, D H Choi, R Pow-754 ell, T Ewalds, P Georgiev, et al. Grandmaster level in star-755 craft ii using multi-agent reinforcement learning. Nature, 756 575(7782):350-354, 2019. 757
- [Watkins and Dayan, 1992] C JCH Watkins and P Dayan. O-758 learning. Machine learning, 8:279-292, 1992. 759
- [Whitley et al., 1993] D Whitley, S Dominic, R Das, and 760 C W Anderson. Genetic reinforcement learning for neu-761 rocontrol problems. Machine Learning, 13(2-3):259-284, 762 1993. 763
- [Williams, 1992] R J Williams. Simple statistical gradient-764 following algorithms for connectionist reinforcement 765 learning. Machine learning, 8:229-256, 1992. 766