

# Back to Basics: Benchmarking Canonical Evolution Strategies for Playing Atari

Patryk Chrabaszcz, Ilya Loshchilov, Frank Hutter

University of Freiburg, Freiburg, Germany

{chrabasp,ilya,fh}@cs.uni-freiburg.de

## Abstract

Evolution Strategies (ES) have recently been demonstrated to be a viable alternative to reinforcement learning (RL) algorithms on a set of challenging deep RL problems, including Atari games and MuJoCo humanoid locomotion benchmarks. While the ES algorithms in that work belonged to the specialized class of natural evolution strategies (which resemble approximate gradient RL algorithms, such as REINFORCE), we demonstrate that even a very basic canonical ES algorithm can achieve the same or even better performance. This success of a basic ES algorithm suggests that the state-of-the-art can be advanced further by integrating the many advances made in the field of ES in the last decades.

We also demonstrate qualitatively that ES algorithms have very different performance characteristics than traditional RL algorithms: on some games, they learn to exploit the environment and perform much better while on others they can get stuck in suboptimal local minima. Combining their strengths with those of traditional RL algorithms is therefore likely to lead to new advances in the state of the art.

## 1 Introduction

In machine learning, Evolution Strategies (ES) are mainly used for **direct policy search** in reinforcement learning [Gomez *et al.*, 2008; Heidrich-Meisner and Igel, 2009; Stulp and Sigaud, 2013; Salimans *et al.*, 2017] and **hyperparameter tuning** in supervised learning, e.g., for Support Vector Machines [Glas-machers and Igel, 2008; Igel, 2011] and Deep Neural Networks [Loshchilov and Hutter, 2016].

Recently it has been shown [Salimans *et al.*, 2017] that ES algorithms can be used for tasks which are dominated by deep reinforcement learning (RL) algorithms. Those tasks include learning a policy with **discrete action set to control an agent's behavior in a wide variety of Atari environments**, as well as learning a policy with continuous action space for agents operating in MuJoCo [Todorov *et al.*, 2012] environments. ES algorithms offer a set of attractive **advantages** when compared to deep RL algorithms:

- They are **highly parallelizable**, since the amount of information that has to be exchanged between workers does not depend on the network size.
- Depending on the problem, they can offer **better exploration**, and as a result different training runs can converge to qualitatively different solutions.
- They are **not sensitive to the distribution of rewards** and do not require careful tuning of discount factors while still facilitating long-term foresight more than traditional discount-based RL algorithms.
- They can be used for the optimization of non-differentiable policy functions.

In this work, we go one step further than Salimans *et al.* [2017] and study the applicability of **even simpler ES** algorithm to the task of learning a policy network for playing Atari games. Salimans *et al.* [2017] used a specialized ES algorithm that belongs to the class of Natural Evolution Strategies (NES) [Wierstra *et al.*, 2008], which computes approximate gradients similar to the REINFORCE algorithm [Williams, 1992]. Here, we demonstrate that very comparable results can already be achieved with a **simpler very basic Canonical ES algorithm from the 1970s** [Rechenberg, 1973; Rudolph, 1997].

Our **contributions** in this work are as follows:

- We demonstrate that even a **very basic** Canonical ES algorithm is able to **match** (and sometimes supersede) the performance of the Natural Evolution Strategy used by Salimans *et al.* [2017] for playing Atari games.
- We demonstrate that after 5 hours of training, Canonical ES is able to find novel solutions that **exploit the game design** and even find bugs in one game that allow them to achieve unprecedented high scores.
- We experimentally study the performance characteristics of both ES algorithms, demonstrating that (1) individual runs have **high variance** in performance and that (2) **longer runs** (5h instead of 1h) lead to significant further performance improvements.
- By demonstrating that Canonical ES is a **competitive alternative** to traditional RL algorithms and the specialized ES algorithms tested so far on the Atari domain we set a benchmark for **future** work on modern ES variants that are directly based on the canonical version.

## 2 Background

In this section, we discuss background on RL for playing Atari and on the previously-introduced NES method.

### 2.1 Reinforcement Learning for Playing Atari

In the Atari task of the OpenAI gym environment [Brockman *et al.*, 2016], the agent needs to learn to maximize its cumulative reward solely by interacting with the environment (i.e., playing the game). Inputs to the agent include the raw pixels displayed on the Atari console as well as the reward signal; its actions correspond to joystick movements to execute.

Recent developments in the field of deep RL have made it possible to address challenging problems that require processing of high dimensional inputs, such as the raw images in this Atari domain, which can be addressed by deep convolutional neural networks. This approach was popularized by Google DeepMind’s Nature paper on the deep Q network (DQN) [Mnih *et al.*, 2015], a Q-learning method that estimates the utility of an action given a current state by means of a deep neural network. Given this network for approximating the Q function, in any state  $s$ , DQN’s policy then simply selects the action  $a$  with the largest predicted Q value  $Q(s, a)$ .

While DQN requires this maximization over the action space, policy gradient algorithms directly parameterize a policy networks that maps a state to a probability distribution over actions. Policy gradient algorithms, such as the Asynchronous Advantage Actor Critic (A3C) [Mnih *et al.*, 2016], directly optimize this policy network.

**State representation.** In Atari games it is important to model the state to include information from previous frames that will influence an agent’s performance. In this work we use the following standard preprocessing pipeline [Mnih *et al.*, 2015] with an implementation provided by OpenAI [Dhariwal *et al.*, 2017]. First, we apply preprocessing to the screen frames to reduce the dimensionality and remove artifacts related to the technical limitations of the Atari game console (flickering). Specifically, we apply a pixel-wise max operation on the current frame and the one preceding it. Next, we convert the output from this operation into a grayscale image, resize it and crop it to 84x84 pixels. At the end of this pipeline, we stack together the result of the 4 last frames produced this way to construct a 84x84x4 state tensor. Also following common practice, to speed up policy evaluation (and thus reduce training time), instead of collecting every frame and making a decision at every step, we collect every 4th frame (3rd frame for SpaceInvaders to make the laser visible [Mnih *et al.*, 2013]) and apply the same action for all frames in between. Figure 1 visualizes the full preprocessing pipeline. Figure 2 shows an example of the state representation for 5 different games.

### 2.2 Natural Evolution for Playing Atari

Salimans *et al.* [2017] recently demonstrated that an ES algorithm from the specialized class of Natural Evolution Strategies (NES; Wierstra *et al.* [2008]) can be used to successfully train policy networks in a set of RL benchmark environments (Atari, MuJoCo) and compete with state-of-the-art RL algorithms. Algorithm 1 describes their approach on a high level. In a nutshell, it evolves a distribution over policy networks

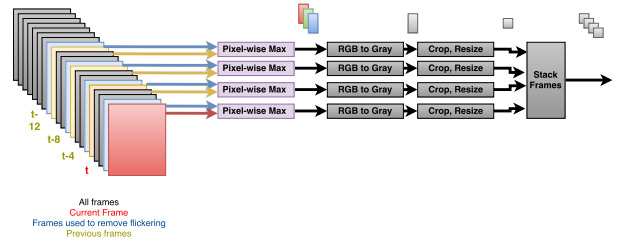


Figure 1: Preprocessing pipeline. Take every 4th frame, apply max operation to remove screen flickering, convert to grayscale, resize/crop, stack 4 last frames.

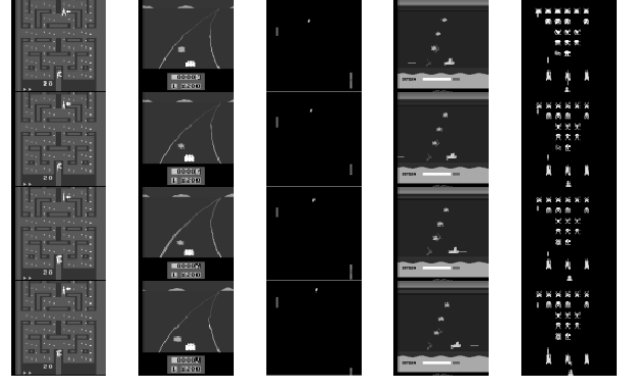


Figure 2: State representation (84x84x4 tensor) for 5 different games: Alien, Enduro, Pong, Seaquest and SpaceInvaders. Channels are shown on top of each other for better visualization.

over time by evaluating a population of  $\lambda$  different networks in each iteration, starting from initial policy parameter vector  $\theta_0$ . At each iteration  $t$ , the algorithm evaluates the game scores  $F(\cdot)$  of  $\lambda$  different policy parameter vectors centered around  $\theta_t$  (lines 3-5) to estimate a gradient signal, using mirrored sampling [Brockhoff *et al.*, 2010] to reduce the variance of this estimate. Since the  $\lambda$  game evaluations are independent of each other, ES can make very efficient use of parallel compute resources. The resulting  $\lambda$  game scores are then ranked (line 6), making the algorithm invariant to their scale; as noted by the authors, this approach (called fitness shaping [Wierstra *et al.*, 2014] but used in all ESs since the 1970s) decreases the probability of falling into local optima early and lowers the influence of outliers. Based on these  $\lambda$  ranks of local steps around  $\theta_t$ , the algorithm approximates a gradient  $g$  (line 7) and uses this with a modern version of gradient descent (Adam [Kingma and Ba, 2014] or SGD with momentum) with weight decay to compute a robust parameter update step (line 8) in order to move  $\theta_t$  towards the parameter vectors that achieved higher scores.

We note that the computation of the approximate gradient  $g$  in line 7 follows the same approach as the well-known policy gradient algorithm REINFORCE. This can be shown as follows. Denoting the distribution from which we draw policy network parameters  $\theta$  as  $p_\psi$ , the gradient of the expected reward  $F(\theta)$  with respect to  $\psi$  is:

$$\nabla_\psi \mathbb{E}_{\theta \sim p_\psi} \{F(\theta)\} = \mathbb{E}_{\theta \sim p_\psi} \{F(\theta) \nabla_\psi \log p_\psi(\theta)\}. \quad (1)$$

Because  $p_\psi$  is chosen as an isotropic Gaussian distribution

---

**Algorithm 1: OpenAI ES**

---

**Input:** $optimizer$  - Optimizer function $\sigma$  - Mutation step-size $\lambda$  - Population size $\theta_0$  - Initial policy parameters $F$  - Policy evaluation function

```

1 for  $t = 0, 1, \dots$  do
2   for  $i = 1, 2, \dots, \frac{\lambda}{2}$  do
3     Sample noise vector:  $\epsilon_i \sim \mathcal{N}(0, I)$ 
4     Evaluate score in the game:  $s_i^+ \leftarrow F(\theta_t + \sigma * \epsilon_i)$ 
5     Evaluate score in the game:  $s_i^- \leftarrow F(\theta_t - \sigma * \epsilon_i)$ 
6     Compute normalized ranks:  $r = ranks(s), r_i \in [0, 1]$ 
7     Estimate gradient:  $g \leftarrow \frac{1}{\sigma * \lambda} \sum_{i=1}^{\lambda} (r_i * \epsilon_i)$ 
8     Update policy network:  $\theta_{t+1} \leftarrow \theta_t + optimizer(g)$ 

```

---

with mean  $\theta_t$  and fixed standard deviation (mutation step-size)  $\sigma$ , the only parameter of  $p_\psi$  is  $\theta_t$  and we have:

$$\nabla_{\psi} \log p_{\psi}(\theta) = \nabla_{\theta_t} \log \frac{1}{\sigma \sqrt{2\pi}} e^{-(\theta - \theta_t)^2 / 2\sigma^2} = \frac{\theta - \theta_t}{\sigma^2} \quad (2)$$

and therefore the following identity holds:

$$\nabla_{\psi} \mathbb{E}_{\theta \sim p_{\psi}} \{F(\theta)\} = \mathbb{E}_{\theta \sim p_{\psi}} \left\{ F(\theta) * \frac{\theta - \theta_t}{\sigma^2} \right\} \quad (3)$$

$$\approx \frac{1}{\sigma * \lambda} \sum_{i=1}^{\lambda} F(\theta^{(i)}) * \frac{(\theta^{(i)} - \theta_t)}{\sigma}, \quad (4)$$

where the last step is simply an approximation by  $\lambda$  samples  $\theta^{(i)} \sim p_{\psi}$ . Equation 4 is exactly as in line 7 of the algorithm except that the raw game scores  $F(\theta^{(i)})$  are replaced with their ranks  $r_i$  due to fitness shaping.

Salimans *et al.* [2017] also made two further contributions to stabilize the training and improve performance. Firstly, they introduced a novel parallelization technique (which uses a noise table to reduce communication cost in order to scale to a large number of  $\lambda$  parallel workers) and used virtual batch normalization [Salimans *et al.*, 2016] to make the network output more sensitive to the noise in the parameter space.

### 3 Canonical ES

While the specialized ES algorithm proposed by OpenAI is equivalent to a **policy gradient algorithm**, in this work we consider a **very basic canonical ES** algorithm that belongs to the prominent family of  $(\mu, \lambda) - ES$  optimization algorithms. Algorithm 2 illustrates this simple approach. Starting from a random parameter vector  $\theta_0$ , in each iteration  $t$  we generate an **offspring** population of size  $\lambda$ . For each element of the population, we add sampled mutation noise  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$  to the current parameter vector  $\theta_t$  (line 3) and evaluate the game score of the resulting vector by one episode rollout (line 4). We then pick the top  $\mu$  parameter vectors according to the collected score and form a new parameter vector  $\theta_{t+1}$  as their weighted mean (lines 5-6).

This algorithm is very **basic in its setting**: we do not use mirrored sampling, we do not decay the parameters, we do

not use any advanced optimizer. The standard weights used to compute the weighted mean of the top  $\mu$  solutions fulfill a similar function to the fitness shaping implemented in OpenAI ES. The new elements introduced by Salimans *et al.* [2017] that we *do* use are **virtual batch normalization** (which is a component of the game evaluations  $F$  and not really of the ES algorithm itself) and the efficient **parallelization of ES** using a random noise table.

We initially implemented the Cumulative Step-size  $\sigma$  Adaptation (CSA) procedure [Hansen and Ostermeier, 1996], which is standard in canonical ES algorithms. However, due to the high time cost of game evaluations, during our time-limited training, we are only able to perform up to thousands update iterations. Since the dimensionality of the parameter vector is relatively large (1.7M), this results in only a negligible change of  $\sigma$  during the training. Therefore, effectively, our algorithm used a fixed step-size and thus we removed step-size adaptation from the description from Algorithm 2 making it even somewhat simpler than typical ES. We employed weighted recombination [Rudolph, 1997] and weights  $w$  as in CSA-ES.

---

**Algorithm 2: Canonical ES Algorithm**

---

**Input:** $\sigma$  - Mutation step-size $\theta_0$  - Initial policy parameters $F$  - Policy evaluation function $\lambda$  - Offspring population size $\mu$  - Parent population size**Initialize :**

$$w_i = \frac{\log(\mu + 0.5) - \log(i)}{\sum_{j=1}^{\mu} \log(\mu + 0.5) - \log(j)}$$

```

1 for  $t = 0, 1, \dots$  do
2   for  $i = 1 \dots \lambda$  do
3     Sample noise:  $\epsilon_i \sim \mathcal{N}(0, I)$ 
4     Evaluate score in the game:  $s_i \leftarrow F(\theta_t + \sigma * \epsilon_i)$ 
5     Sort  $(\epsilon_1, \dots, \epsilon_{\lambda})$  according to  $s$  ( $\epsilon_i$  with best  $s_i$  first)
6     Update policy:  $\theta_{t+1} \leftarrow \theta_t + \sigma * \sum_{j=1}^{\mu} w_j * \epsilon_j$ 
7     Optionally, update step size  $\sigma$  (see text)

```

---

## 4 Experiments

In our experiments, we evaluate the performance of the Canonical ES on a subset of 8 Atari games available in OpenAI Gym [Brockman *et al.*, 2016]. We selected these games to represent different levels of difficulty, ranging from simple ones like Pong and Breakout to complex games like Qbert and Alien. We make our implementation of the Canonical ES algorithm available online at [https://github.com/PatrykChrabaszcz/Canonical\\_ES\\_Atari](https://github.com/PatrykChrabaszcz/Canonical_ES_Atari).

We compare our results against those obtained with the ES algorithm proposed by OpenAI [Salimans *et al.*, 2017]. Since no implementation of that algorithm is publicly available for Atari games, we re-implemented it with help from OpenAI<sup>1</sup> and the results of our implementation (which we refer to as “OpenAI ES (our)”) roughly match those reported by Salimans *et al.* [2017] (see Table 1).

<sup>1</sup>We thank Tim Salimans for his helpful email support.

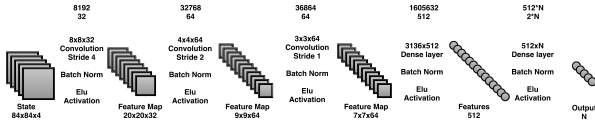


Figure 3: Neural network architecture. Numbers on top show number of parameters in each layer (kernel parameters and batch norm parameters). Each batch norm layer has a trainable shift parameter  $\beta$ ; the last batch norm has an additional trainable scale parameter  $\alpha$ .

**Network Architecture.** We use the same network structure as the original DQN work [Mnih *et al.*, 2015], only changing the activation function from ReLU to ELU [Clevert *et al.*, 2015] and adding batch normalization layers [Ioffe and Szegedy, 2015]. The network as presented in Figure 3 has approximately 1.7M parameters. We initialize network weights using samples from a normal distribution  $\mathcal{N}(\mu = 0, \sigma = 0.05)$ .

**Virtual Batch Normalization.** Following Salimans *et al.* [2017], we use virtual batch normalization [Salimans *et al.*, 2016]. In order to collect the reference batch, at the beginning of the training we play the game using random actions. In each step, we save the corresponding state with the probability  $p(\text{save}) = 1\%$  and stop when 128 samples have been collected.

**Training.** For each game and each ES variant we tested, we performed 3 training runs, each on 400 CPUs with a time budget of 10 hours. Every worker (CPU) evaluates 2 offspring solutions, meaning that our setting is roughly the same as training for 5 hours with full parallelization (800 CPUs); therefore, we label this setting as “5 hours”. In addition, we save the solution proposed after 2 hours of training (equivalent to 1 hour with full parallelization) or after 1 billion training frames (whatever comes first) to allow for a fair comparison with results reported by Salimans *et al.* [2017]; we label this setting as “1 hour”. During training, one CPU is reserved to evaluate the performance of the solution proposed in the current iteration; hence, the offspring population size in our experiments is  $\lambda = 798$ . In each decision step, the agent passes its current environment state through the network and performs an action that corresponds to the output with the highest value. We limit episodes to have a maximum length of 25k steps; we do not adjust this value during training.

An episode includes multiple lives, and we do not terminate an episode after the agent dies the first time in order to allow the learning of strategies that span across multiple lives. We start each episode with up to 30 initial random no-op actions.

**Results.** First, we studied the importance of the parent population size  $\mu$ . This hyperparameter is known to often be important for ES algorithms and we found this to also hold here. We measured performance for  $\mu \in \{10, 20, 50, 100, 200, 400\}$  and observed different optimal values of  $\mu$  for different games (Figure 4). For the subsequent analyses we fixed  $\mu = 50$  for all games.

In Table 1, for each game and for both Canonical ES and OpenAI ES, we report the results of our 3 training runs; for each of these runs, we evaluated the final policy found using 30 rollouts. We ordered the 3 runs according to their mean

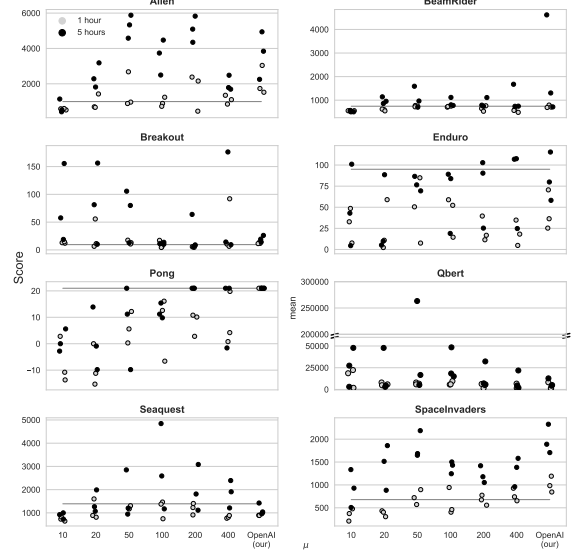


Figure 4: Final evaluation scores (mean across 30 evaluation rollouts with random initial no-ops) for Canonical ES ( $\mu \in \{10, 20, 50, 100, 200, 400\}$ ) and for our implementation of OpenAI ES. For each setting we use population size  $\lambda = 798$  and report the performance from 3 separate training runs. We report evaluation scores after 1 hour and after 5 hours of training. The horizontal line indicates the results reported by Salimans *et al.* [2017].

score and compared the results row-wise. We ran a Mann-Whitney U test [Mann and Whitney, 1947] to check if the distributions of the 30 rewards significantly differed between the two ES variants. After 5 hours of training, Canonical ES performed significantly better than OpenAI ES in 9 out of 24 different runs ( $p < 0.05$ ) and worse in 7 (with 8 ties); this shows that our simple Canonical ES is competitive with the specialized OpenAI ES algorithm on complex benchmark problems from the Atari domain. Additionally, we tested whether the algorithms still made significant progress between 1 and 5 hours of training; indeed, the performance of our Canonical ES algorithm improved significantly in 16 of 24 cases, demonstrating that it often could make effective use of additional training time. However, qualitatively, we observed that the performance of both algorithms tends to plateau in locally optimal solutions for extended periods of time (see Figure 5) and that they often find solutions that are not robust to the noise in the domain; i.e., there is a high variance in the points scored with the same trained policy network across initial environment conditions.

## 5 Qualitative analysis

Visual inspection and comparison between solutions found by reinforcement learning algorithms and solutions found by ES algorithms shows some significant differences. In this section we describe interesting agent behaviors on different games; a video with evaluation runs on all games is available on-line at <https://www.youtube.com/watch?v=0wDzPBiURSI>.

First, we study two games in which most of the ES runs converged to similar sub-optimal solutions: Seaquest and En-



	OpenAI ES 1 hour	OpenAI ES (our) 1 hour	Canonical ES 1 hour	OpenAI ES (our) 5 hours	Canonical ES 5 hours
Alien		<b>3040 <math>\pm</math> 276.8</b>	2679.3 $\pm$ 1477.3	4940 $\pm$ 0	<b>5878.7 <math>\pm</math> 1724.7</b>
Alien	994	<b>1733.7 <math>\pm</math> 493.2</b>	965.3 $\pm$ 229.8	3843.3 $\pm$ 228.7	<b>5331.3 <math>\pm</math> 990.1</b>
Alien		<b>1522.3 <math>\pm</math> 790.3</b>	885 $\pm$ 469.1	2253 $\pm$ 769.4	<b>4581.3 <math>\pm</math> 299.1</b>
BeamRider		<b>792.3 <math>\pm</math> 146.6</b>	774.5 $\pm$ 202.7	<b>4617.1 <math>\pm</math> 1173.3</b>	1591.3 $\pm$ 575.5
BeamRider	744	708.3 $\pm$ 194.7	<b>746.9 <math>\pm</math> 197.8</b>	<b>1305.9 <math>\pm</math> 450.4</b>	965.3 $\pm$ 441.4
BeamRider		690.7 $\pm$ 87.7	<b>719.6 <math>\pm</math> 197.4</b>	<b>714.3 <math>\pm</math> 189.9</b>	703.5 $\pm$ 159.8
Breakout		14.3 $\pm$ 6.5	<b>17.5 <math>\pm</math> 19.4</b>	26.1 $\pm$ 5.8	<b>105.7 <math>\pm</math> 158</b>
Breakout	9.5	11.8 $\pm$ 3.3	<b>13 <math>\pm</math> 17.1</b>	19.4 $\pm$ 6.6	<b>80 <math>\pm</math> 143.4</b>
Breakout		<b>11.4 <math>\pm</math> 3.6</b>	10.7 $\pm$ 15.1	<b>14.2 <math>\pm</math> 2.7</b>	12.7 $\pm$ 17.7
Enduro		70.6 $\pm$ 17.2	<b>84.9 <math>\pm</math> 22.3</b>	<b>115.4 <math>\pm</math> 16.6</b>	86.6 $\pm$ 19.1
Enduro	95	36.4 $\pm$ 12.4	<b>50.5 <math>\pm</math> 15.3</b>	<b>79.9 <math>\pm</math> 18</b>	76.5 $\pm$ 17.7
Enduro		<b>25.3 <math>\pm</math> 9.6</b>	7.6 $\pm$ 5.1	58.2 $\pm$ 10.5	<b>69.4 <math>\pm</math> 32.8</b>
Pong		<b>21.0 <math>\pm</math> 0.0</b>	12.2 $\pm$ 16.6	<b>21.0 <math>\pm</math> 0.0</b>	<b>21.0 <math>\pm</math> 0.0</b>
Pong	21	<b>21.0 <math>\pm</math> 0.0</b>	5.6 $\pm$ 20.2	<b>21 <math>\pm</math> 0</b>	11.2 $\pm$ 17.8
Pong		<b>21.0 <math>\pm</math> 0.0</b>	0.3 $\pm$ 20.7	<b>21 <math>\pm</math> 0</b>	-9.8 $\pm$ 18.6
Qbert		<b>8275 <math>\pm</math> 0</b>	8000 $\pm$ 0	12775 $\pm$ 0	<b>263242 <math>\pm</math> 433050</b>
Qbert	147.5	1400 $\pm$ 0	<b>6625 <math>\pm</math> 0</b>	5075 $\pm$ 0	<b>16673.3 <math>\pm</math> 6.2</b>
Qbert		1250 $\pm$ 0	<b>5850 <math>\pm</math> 0</b>	4300 $\pm$ 0	<b>5136.7 <math>\pm</math> 4093.9</b>
Seaquest		1006 $\pm$ 20.1	<b>1306.7 <math>\pm</math> 262.7</b>	1424 $\pm$ 26.5	<b>2849.7 <math>\pm</math> 599.4</b>
Seaquest	1390	898 $\pm$ 31.6	<b>1188 <math>\pm</math> 24</b>	1040 $\pm$ 0	<b>1202.7 <math>\pm</math> 27.2</b>
Seaquest		887.3 $\pm$ 20.3	<b>1170.7 <math>\pm</math> 23.5</b>	<b>960 <math>\pm</math> 0</b>	946.7 $\pm$ 275.1
SpaceInvaders		<b>1191.3 <math>\pm</math> 84.6</b>	896.7 $\pm$ 123	<b>2326.5 <math>\pm</math> 547.6</b>	2186 $\pm$ 1278.8
SpaceInvaders	678.5	<b>983.7 <math>\pm</math> 158.5</b>	721.5 $\pm$ 115	<b>1889.3 <math>\pm</math> 294.3</b>	1685 $\pm$ 648.6
SpaceInvaders		<b>845.3 <math>\pm</math> 69.7</b>	571.3 $\pm$ 98.8	<b>1706.5 <math>\pm</math> 118.3</b>	1648.3 $\pm$ 294.5

Table 1: Evaluation scores (mean over 30 evaluation runs with up to 30 initial no-ops) for different training times and algorithms. For each training time limit we compare results from OpenAI ES(our) and Canonical ES  $\mu = 50$ . For each setting we performed 3 training runs, ordered the results for each game and compared them row by row, boldfacing the better score. Results for which the difference is significant across the 30 evaluation runs based on a Mann-Whitney U test ( $p < 0.05$ ) are marked in blue.

duro. In Seaquest, the agent dives to the bottom of the sea and starts to shoot left and right, occasionally hitting an enemy and scoring points (Figure 6). However, it is not able to detect the lack of oxygen and quickly dies. In Enduro, the agent steers the car to keep it in the middle of the road and accelerate, but from time to time it bounces back after hitting a rival car in front of it. Both solutions are easy-to-reach local optima and do not require developing any complex behavior; since the corresponding scores achieved with these policies are much higher than those of random policies we believe that it is hard for ES to escape from these local attractor states in policy space.

We next study the game Qbert, in which Canonical ES found two particularly interesting solutions. In the first case (<https://www.youtube.com/watch?v=-p7VhdTXA0k>), the agent gathers some points at the beginning of the game and then stops showing interest in completing the level. Instead, it starts to bait an enemy that follows it to kill itself. Specifically, the agent learns that it can jump off the platform when the enemy is right next to it, because the enemy will follow: although the agent loses a life, killing the enemy yields enough points to gain an extra life again (Figure 7). The agent repeats this cycle of suicide and killing the opponent over and over again.

In the second interesting solution (<https://www.youtube.com/watch?v=meE5aaRJ0Zs>), the agent discovers an in-game bug. First, it completes the first level and then starts to jump from platform to platform in what seems

to be a random manner. For a reason unknown to us, the game does not advance to the second round but the platforms start to blink and the agent quickly gains a huge amount of points (close to 1 million for our episode time limit). Interestingly, the policy network is not always able to exploit this in-game bug and 22/30 of the evaluation runs (same network weights but different initial environment conditions) yield a low score.

Breakout seems to be a challenging environment for ES algorithms. Canonical ES only found reasonable solutions for a few settings. The best solution shows a strategy that looks similar to the best strategies obtained by using reinforcement learning algorithms, in which the agent creates a hole on one side of the board and shoots the ball through it to gain many points in a short amount of time (Figure 8). However, even this best solution found by ES is not stable: for different initial environment conditions the agent with the same policy network quickly loses the game with only few points.

In the games SpaceInvaders and Alien we also observed interesting strategies. We do not clip rewards during the training as is sometimes done for reinforcement learning algorithms. Because of that the agent puts more attention to behaviors that result in a higher reward, sometimes even at the cost of the main game objective. In SpaceInvaders, we observe that in the best solution the agent hits the mother-ship that appears periodically on the top of the screen with 100% accuracy. In Alien, the agent focuses on capturing an item that makes it invincible and then goes to the enemy spawn point to collect rewards for eliminating newly appearing enemies. However,

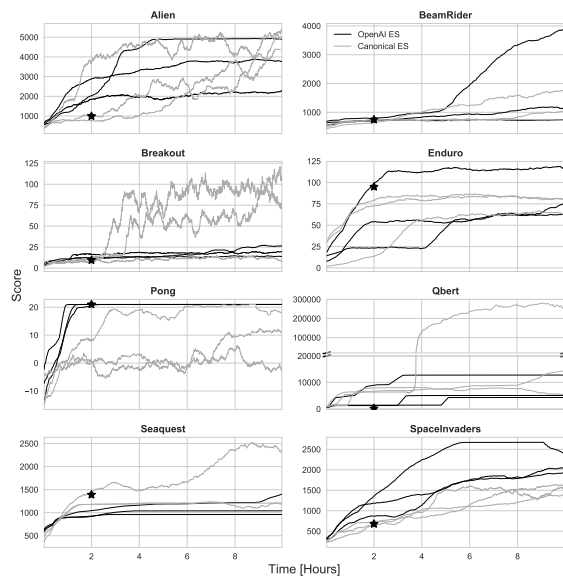


Figure 5: Training curves for CanonicalES ( $\mu = 50$ ) and OpenAI ES (our). At each iteration  $t$  we evaluate currently proposed solution  $\theta_t$  two times using one CPU, because of that values reported in the Table 1 (mean over 30 evaluation runs) might differ. For better plot quality we filter out the noise. We observe that both algorithms often get stuck in local optima.

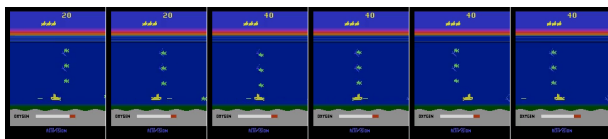


Figure 6: The agent learns to dive to the bottom of the sea and constantly shoot left and right, occasionally scoring points.

the agent is not able to detect when the invincibility period ends.

## 6 Recent related work

Evolutionary strategies, such as the Covariance Matrix Adaptation Evolution Strategy (CMA-ES [Hansen *et al.*, 2003]), are commonly used as a baseline approach in reinforcement learning tasks [Heidrich-Meisner and Igel, 2009; Stulp and Sigaud, 2013; Duan *et al.*, 2016; Li, 2017]. Here, we only discuss the most recent related works, which also followed up on the work by Salimans *et al.* [2017]; in particular, we discuss three related arXiv preprints that scientists at Uber released in the last two months about work concurrent to ours.

Similarly to our work, Such *et al.* [2017] studied the performance of simpler algorithms than OpenAI’s specialized ES variant. They show that genetic algorithms (another broad class of black-box optimization algorithms) can also reach results competitive to OpenAI’s ES variant and other RL algorithms. Additionally, interestingly, the authors show that for some of the games even simple random search can outperform carefully designed RL and ES algorithms.

Lehman *et al.* [2017] argue that comparing ES to finite-difference-based approximation is too simplistic. The main

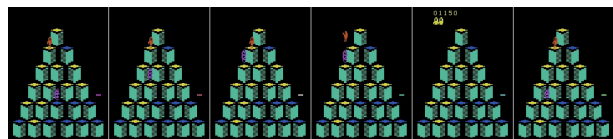


Figure 7: The agent (orange blob in the upper left part of the screen) learns to commit suicide to kill its enemy (purple spring) and collects enough points to get another life. The whole cycle is repeated over and over again.

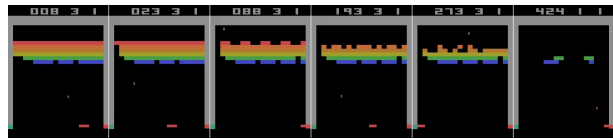


Figure 8: The agent learns to make a hole in the brick wall to collect many points with one ball bounce.

difference comes from the fact that ES tries to optimize the performance of the distribution of solutions rather than a single solution, thus finding solutions that are more robust to the noise in the parameter space. The authors leave open the question whether this robustness in the parameter space also affects the robustness to the noise in the domain. In our experiments we observe that even for the best solutions on some of the games, the learned policy network is not robust against environment noise.

Conti *et al.* [2017] try to address the problems of local minima (which we also observed in games like Seaquest and Enduro) by augmenting ES algorithms with a novelty search (NS) and quality diversity (QD). Their proposed algorithms add an additional criterion to the optimization procedure that encourages exploring qualitatively different solutions during training, thus reducing the risk of getting stuck in a local optimum. The authors also propose to manage a meta-population of diverse solutions allocating resources to train more promising ones. In our experiments we observe that training runs with the same hyperparameters and different initializations often converge to achieve very different scores; managing a meta-population could therefore be an easy way to improve the results and reduce the variance between the runs. Overall, the success of Conti *et al.* [2017] in improving performance with some of these newer methods in ES research strengthens our expectation that a wide range of improvements to the state of the art are possible by integrating the multitude of techniques developed in ES research over the last decades into our canonical ES.

## 7 Conclusion

The recent results provided by Open AI [Salimans *et al.*, 2017] suggest that natural evolution strategies represent a viable alternative to more common approaches used for deep reinforcement learning. In this work, we analyzed whether the results demonstrated in that work are due to the special type of evolution strategy used there. Our results suggest that even a very basic decades-old evolution strategy provides comparable results; thus, more modern evolution strategies should be considered as a potentially competitive approach to modern deep reinforcement learning algorithms. Since evolution

strategies have different strengths and weaknesses than traditional deep reinforcement learning strategies, we also expect rich opportunities for combining the strength of both.

## Acknowledgments

This work has partly been supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme under grant no. 716721. The authors acknowledge support by the state of Baden-Württemberg through bwHPC and the German Research Foundation (DFG) through grant no INST 39/963-1 FUGG.

## References

- [Brockhoff *et al.*, 2010] Dima Brockhoff, Anne Auger, Nikolaus Hansen, Dirk V Arnold, and Tim Hohm. Mirrored sampling and sequential selection for evolution strategies. In *International Conference on Parallel Problem Solving from Nature*, pages 11–21. Springer, 2010.
- [Brockman *et al.*, 2016] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [Clevert *et al.*, 2015] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- [Conti *et al.*, 2017] Edoardo Conti, Vashisht Madhavan, Felipe Petroski Such, Joel Lehman, Kenneth O Stanley, and Jeff Clune. Improving exploration in evolution strategies for deep reinforcement learning via a population of novelty-seeking agents. *arXiv preprint arXiv:1712.06560*, 2017.
- [Dhariwal *et al.*, 2017] Prafulla Dhariwal, Christopher Hesse, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, and Yuhuai Wu. Openai baselines. <https://github.com/openai/baselines>, 2017.
- [Duan *et al.*, 2016] Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *ICML*, pages 1329–1338, 2016.
- [Glasmachers and Igel, 2008] Tobias Glasmachers and Christian Igel. Uncertainty handling in model selection for support vector machines. In *Parallel Problem Solving from Nature (PPSN)*, pages 185–194. Springer, 2008.
- [Gomez *et al.*, 2008] Faustino Gomez, Jürgen Schmidhuber, and Risto Miikkulainen. Accelerated neural evolution through cooperatively coevolved synapses. *Journal of Machine Learning Research*, 9(May):937–965, 2008.
- [Hansen and Ostermeier, 1996] Nikolaus Hansen and Andreas Ostermeier. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation. In *Evolutionary Computation, 1996., Proceedings of IEEE International Conference on*, pages 312–317. IEEE, 1996.
- [Hansen *et al.*, 2003] Nikolaus Hansen, Sibylle D Müller, and Petros Koumoutsakos. Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (cma-es). *Evolutionary computation*, 11(1):1–18, 2003.
- [Heidrich-Meisner and Igel, 2009] Verena Heidrich-Meisner and Christian Igel. Hoeffding and bernstein races for selecting policies in evolutionary direct policy search. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 401–408. ACM, 2009.
- [Igel, 2011] Christian Igel. Evolutionary kernel learning. In *Encyclopedia of Machine Learning*, pages 369–373. Springer, 2011.
- [Ioffe and Szegedy, 2015] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, pages 448–456, 2015.
- [Kingma and Ba, 2014] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [Lehman *et al.*, 2017] Joel Lehman, Jay Chen, Jeff Clune, and Kenneth O Stanley. Es is more than just a traditional finite-difference approximator. *arXiv preprint arXiv:1712.06568*, 2017.
- [Li, 2017] Yuxi Li. Deep reinforcement learning: An overview. *arXiv preprint arXiv:1701.07274*, 2017.
- [Loshchilov and Hutter, 2016] Ilya Loshchilov and Frank Hutter. CMA-ES for Hyperparameter Optimization of Deep Neural Networks. *arXiv preprint arXiv:1604.07269*, 2016.
- [Mann and Whitney, 1947] Henry B Mann and Donald R Whitney. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics*, pages 50–60, 1947.
- [Mnih *et al.*, 2013] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [Mnih *et al.*, 2015] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [Mnih *et al.*, 2016] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *ICML*, pages 1928–1937, 2016.
- [Rechenberg, 1973] Ingo Rechenberg. Evolutionsstrategie—optimierung technischer systeme nach prinzipien der biologischen evolution. 1973.
- [Rudolph, 1997] Günter Rudolph. *Convergence properties of evolutionary algorithms*. Kovac, 1997.
- [Salimans *et al.*, 2016] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pages 2234–2242, 2016.
- [Salimans *et al.*, 2017] Tim Salimans, Jonathan Ho, Xi Chen, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- [Stulp and Sigaud, 2013] Freek Stulp and Olivier Sigaud. Robot skill learning: From reinforcement learning to evolution strategies. *Paladyn, Journal of Behavioral Robotics*, 4(1):49–61, 2013.
- [Such *et al.*, 2017] Felipe Petroski Such, Vashisht Madhavan, Edoardo Conti, Joel Lehman, Kenneth O Stanley, and Jeff Clune. Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning. *arXiv preprint arXiv:1712.06567*, 2017.
- [Todorov *et al.*, 2012] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In

*IROS, 2012 IEEE/RSJ International Conference on*, pages 5026–5033. IEEE, 2012.

[Wierstra *et al.*, 2008] Daan Wierstra, Tom Schaul, Jan Peters, and Juergen Schmidhuber. Natural evolution strategies. In *Evolutionary Computation, 2008. CEC 2008.(IEEE World Congress on Computational Intelligence). IEEE Congress on*, pages 3381–3387. IEEE, 2008.

[Wierstra *et al.*, 2014] Daan Wierstra, Tom Schaul, Tobias Glas-machers, Yi Sun, Jan Peters, and Jürgen Schmidhuber. Natural evolution strategies. *Journal of Machine Learning Research*, 15(1):949–980, 2014.

[Williams, 1992] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. In *Reinforcement Learning*, pages 5–32. Springer, 1992.