

Hybrid Deep Reinforcement Learning: Combine the best of Evolutionary Computation and Policy Gradient Methods

Ziyu Lu

August 20, 2018

1 Introduction

In the recent years, with the development of deep learning and reinforcement learning, deep reinforcement learning (DRL) has become an important research area in the AI research. The main idea of deep reinforcement learning is to apply deep neural network to traditional reinforcement learning problems, mirroring how humans learn to best exploit the environment from exploring and interacting with the environment. Among different DRL algorithms, Policy Gradient methods have been popular among scientists for years, and a lot of variations, including Actor-Critic[1], TRPO[2], PPO[3], have been proposed since the emergence of the vanilla Policy Gradient method, REINFORCE[4]. These methods are all based on the the idea of Markov Decision Process, and the neural networks are optimized through back-propagating gradients. Nevertheless, recent works have shown that a class of black box optimization algorithms inspired by natural evolution and selection, such as Evolution Strategies (ES)[5] and Genetic Algorithms (GA)[6], can serve as a strong alternative to the Policy Gradient methods. Performing direct policy search, these algorithms require no policy differentiation and no back-propagation. Different as they are, the two families of algorithms have their respective advantages and disadvantages. For example, while convergence is usually guaranteed for Policy Gradient methods, there is always concern for the policy getting stuck in local optima; and while ES and GA explore more aggressively in the policy space than Policy Gradient methods, they suffer from poorer sample complexity compared to their gradient-based counterpart on harder environments. In the summer of 2018, under the supervision of Prof. Keith Ross, Dean of Engineering and Computer Science at NYU Shanghai, and Che Wang, Ph.D. candidate in Computer Science at NYU, we have been experimenting with various approaches of combining these two families of algorithms together, hoping to find an approach that can use the advantage of one kind of algorithms to address the problem of the other, or even better, that can combines their advantages together.

2 Combine Evolution Strategy with vanilla Policy Gradient

Algorithm 1 ES/PG, with adaptive weight

Input: ES learning rate α , noise standard deviation σ , initial policy parameters θ_0 , initial policy gradient weight w_{pg}

for $t = 0, 1, 2, \dots$ **do**

- Sample $\epsilon_1, \dots, \epsilon_n \sim N(0, I\sigma)$
- Compute returns $F_i = F(\theta_t + \sigma\epsilon_i)$ for $i = 1, \dots, n$
- Compute policy gradient g of θ_t
- $\theta_{t+1} \leftarrow \theta_t + w_{pg}g + (1 - w_{pg})\alpha \frac{1}{n\sigma} \sum_{i=1}^n F_i \epsilon_i$
- $w_{pg} \leftarrow 10 * w_{pg}$

The idea of adaptive weights comes from the observation that this combination generates decent results only when the weight for policy gradient is sufficiently small. Experiment results with some fixed weights on the InvertedPendulum-v2 environment also indicates that smaller weights for policy gradient tend to generate better results. Those experiment results can be seen in the figure 1.

Note: “no pg” refers to vanilla ES; w_{pg} refers to the weight for policy gradient in each experiment.

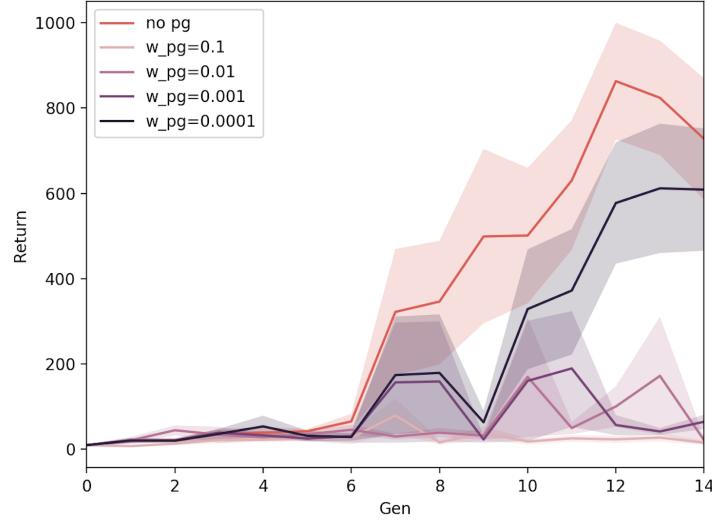


Figure 1

In addition, there is another observation that the performance of Evolution Strategy can still be volatile after the optimum is reached. This observation probably comes from the fact that Evolution Strategy doesn't guarantee convergence. Since Policy Gradient

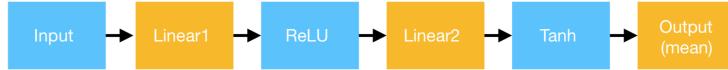
methods usually have a more stable performance improvement and the gradient of a policy tends to zero as the policy reaches the optimum, we hope that by increasing the weight for policy gradient as more generations have been generated, the policy can be more stable after it reaches the optimum.

Since we want the weight for policy gradient to be small enough at the beginning and gradually get larger, and the sum of the weight for the gradient and the weight for Evolution Strategy update to be one, in each experiment, we started with the weight for policy gradient set to be less than or equal to $10^{-\text{total gen}}$, and the weight for policy gradient becomes 10 times larger after each generation.

2.1 Experiment results on InvertedPendulum environment

Below are some results we got on the InvertedPendulum-v2 environment as well as the settings we used:

- Neural network architecture:



- Standard deviation for sampling actions: 0.1
- Number of hidden units: 4
- Population size: 100
- α : 5
- σ : 1
- Policy gradient batch size: 1000
- Results were averaged over 7 random seeds (1,2,3,4,5,6,7)

2.1.1 ES/PG, with adaptive weight

Note: “no pg” refers to vanilla ES; w_{pg} refers to the initial weight for policy gradient in each experiment; the learning curves “no pg” and “ $w_{\text{pg}}=1e-25$ ” overlap each other.

In figure 2, it is clear that experiments with the weight for policy gradient initialized to $1e-20$ or $1e-25$ give decent results. Meanwhile, the performance of Algorithm 1 is not so satisfactory if the initial weight for policy gradient is set to $1e-15$. In particular, its performance suffers from significant decrease in the last few generations, as the weight for policy gradient increases to 0.01 or 0.1.

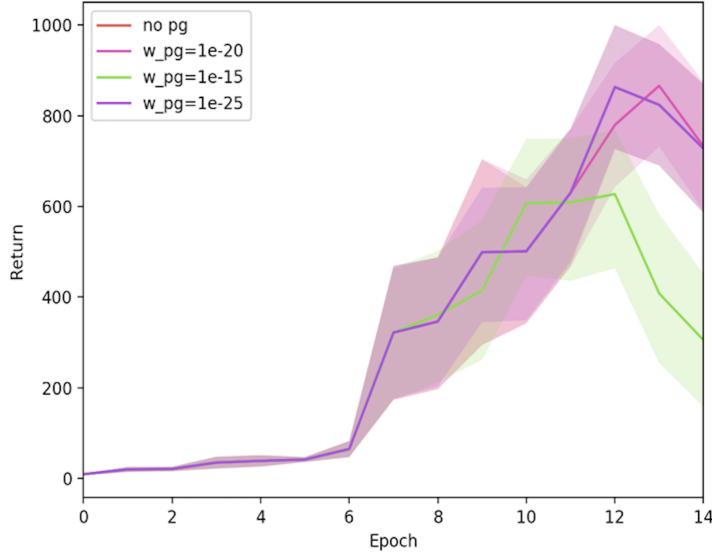


Figure 2

2.1.2 ES/PG, with ES weight decay

To see the impact of the gradients, we also experimented with decreasing the weight for Evolution Strategy update as in Algorithm 1, but without policy gradient computed. The results are shown in figure 3.

Note: The learning curves for “no pg”, “w_pg=1e-25”, as well as “w_es decay=1e-20” overlap each other; “w_pg decay=1e-15” means that in the nth generation ($n \geq 0$) the weight for evolution strategy is equal to $1 - (1e - 15) * (10^n)$.

Among these methods, it seems that “w_es decay=1e-15” gives better result, while “w_pg=1e-20”, “w_pg=1e-25”, and “w_es decay=1e-20” give decent results as well.

2.1.3 ES/PG, with fixed weight

We also experimented with small fixed weights for policy gradient to see whether fixed weights can give decent results as well if small enough. The results are shown in figure 4.

Among the weights tested, 5e-6 achieves better result than vanilla Evolution Strategy in the first 7 generations, but none of them has performance comparable to that of vanilla Evolution Strategy by the end of the experiment.

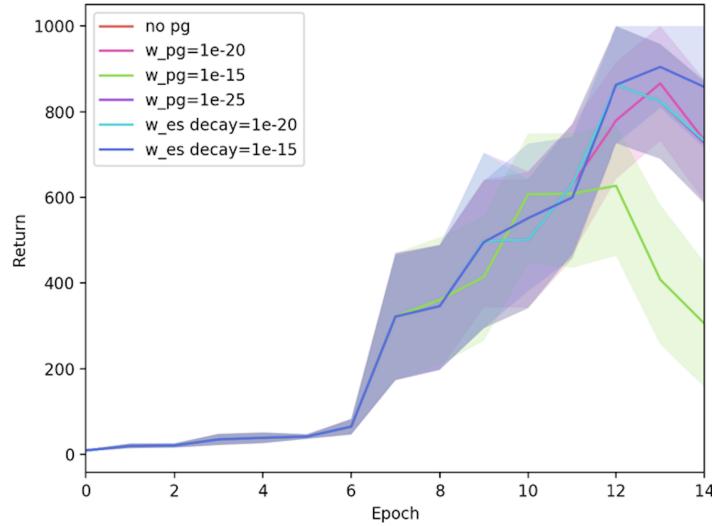


Figure 3

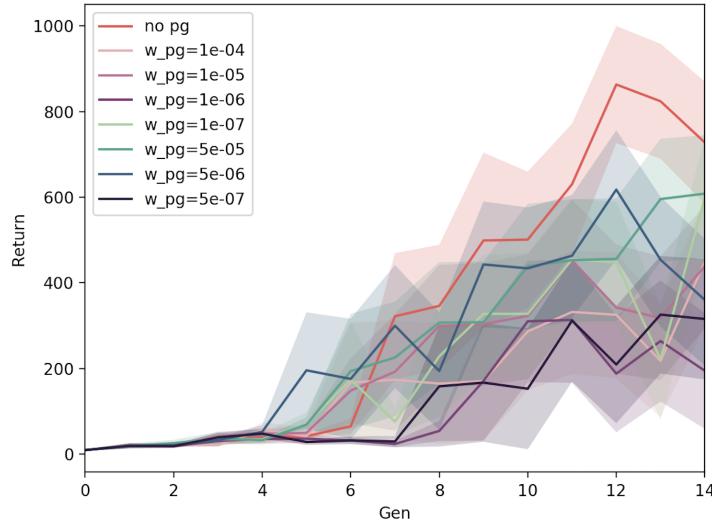


Figure 4

2.1.4 ES/PG, with adaptive weight and fixed weight

Figure 5 compares some results generated with fixed weights and adaptive weights for policy gradient, as well as the result of vanilla Evolution Strategy.

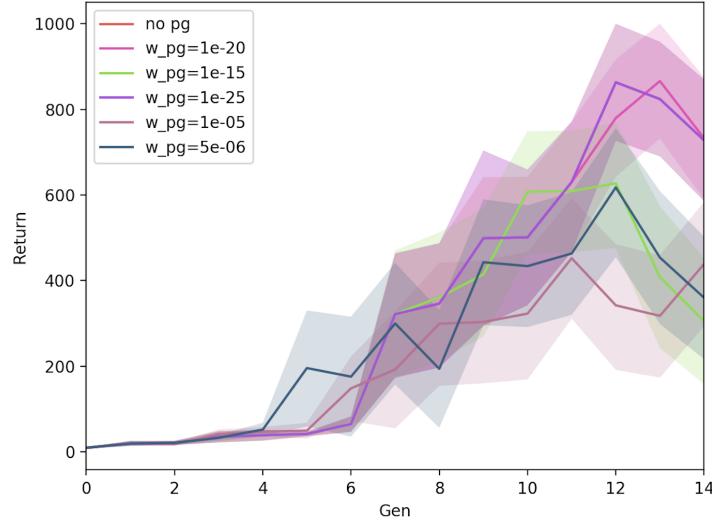


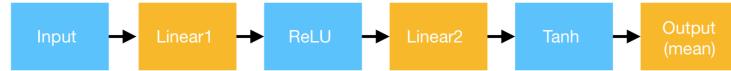
Figure 5

It can be seen from figure 5 that although fixed weight 5e-6 outperforms all other configurations during the first 7 generations, only adaptive policy gradient weights can achieve performances comparable to that of vanilla Evolution Strategy by the end of experiments.

2.2 Experiment results on HalfCheetah environment

Below are some results as well as the settings on the HalfCheetah-v2 environment:

- Neural network architecture:



- Standard deviation for sampling actions: 0.05
- Number of hidden units: 32
- Population size: 100
- α : 0.1
- σ : 0.05
- Results were averaged over 5 random seeds (1,2,3,4,5)

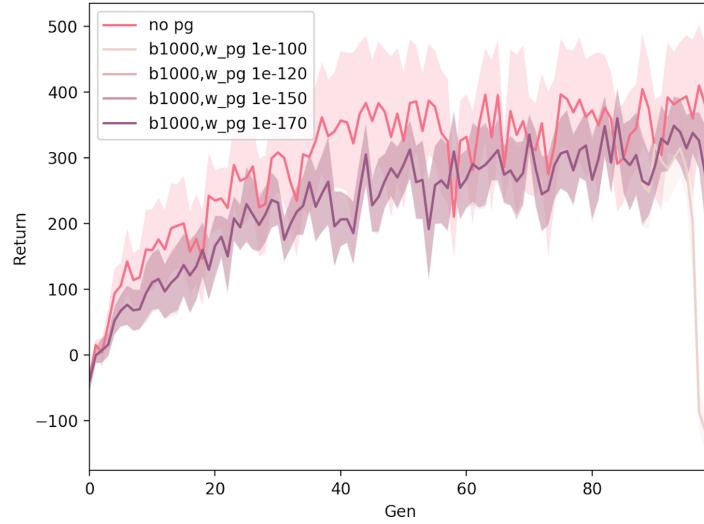


Figure 6: ES&PG, 100 generations, maximum path length 200, policy gradient batch size 1000

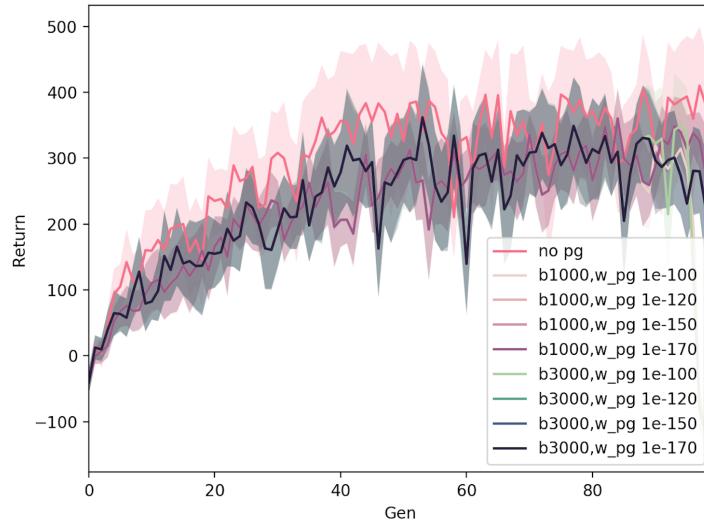


Figure 7: ES&PG, 100 generations, maximum path length 200, policy gradient batch size 1000&3000

Figures 6 and 7 compare performances of vanilla ES and Algorithm 1 with different policy gradient batch sizes.

Among the weights tested, it seems that vanilla ES still gives the best result. And

from figure 7 it seems that larger batch size doesn't make a significant improvement on the performances of Algorithm 1. We didn't test weight decay on this environment, because as the number of generations is large in this case, weight decay for Evolution Strategy update won't be effective until the last few generations.

3 Combine Evolution Strategy with advanced Policy Gradient methods

As vanilla Evolution Strategy always outperforms vanilla Policy Gradient as well as Algorithm 1, it seems to us that this naive combination of ES update and Policy Gradient update is not very promising. We then proposed a new algorithm, which performs one update, either ES or PG, at a time, instead of combining two updates together.

Algorithm 2 ES/PG, with learning goal

Input: ES learning rate α_1 , PG learning rate α_2 , noise standard deviation σ , initial policy parameters θ_0 , learning goal G

for $t = 0, 1, 2, \dots$ **do**

- Compute return $F(\theta_t)$
- if** $F(\theta_t) < 0.8 * G$ **then**
- Sample $\epsilon_1, \dots, \epsilon_n \sim N(0, I\sigma)$
- Compute returns $F_i = F(\theta_t + \sigma\epsilon_i)$ for $i = 1, \dots, n$
- $\theta_{t+1} \leftarrow \theta_t + \alpha_1 \frac{1}{n\sigma} \sum_{i=1}^n F_i \epsilon_i$
- else**
- Compute policy gradient $\nabla_\theta J(\theta)$ of θ_t
- $\theta_{t+1} \leftarrow \theta_t + \alpha_2 \nabla_\theta J(\theta)$

Algorithm 2 will perform Evolution Strategy update when the return of the center network is below the threshold (80% learning goal), typically at the beginning of experiments. By doing so, we are able to take advantage of the observation that Evolution Strategy explores more aggressively and learns faster (especially at the beginning of the experiments) than Policy Gradient methods. Once the return of the center network reaches the threshold, we will follow the policy gradient update in that generation. Since the policy seems to be close to an optimum, it is reasonable to reduce exploration and approaches the optimum in a more steady fashion.

In order to avoid performance collapse after switching to policy gradient update from Evolution Strategy update, we need to find a policy gradient method of which the performance can approximately match that of vanilla Evolution Strategy. Since the performance of vanilla Policy Gradient is not able to match that of vanilla Evolution Strategy, we implemented Algorithm 2 with two more advanced policy gradient methods, the Actor-Critic algorithm and the PPO algorithm.

3.1 Combine Evolution Strategy with Actor-Critic

When implementing Algorithm 2 with Actor-Critic algorithm, we experimented with updating baseline multiple times using mini-batch gradient descent, as well as using different learning rates for the Actor and Critic networks.

3.1.1 Multiple updates on baseline network

As the prediction accuracy of the baseline network is critical to the performance of the Actor-Critic algorithm, we update the baseline network not only in each policy gradient step, but also in each Evolution Strategy generation, using data generated from fitness evaluation. We also tried mini-batch gradient descent by dividing the training data into several mini-batches and then running one epoch.

Figures 8-13 compare the performances of Actor-Critic algorithm with different learning rates as well as different mini-batch sizes on the HalfCheetah-v2 environment. Neural network architecture and hyper-parameter setting are consistent in all these experiments, and are listed below:

Actor:



- Standard deviation for sampling actions: 0.05
- Number of hidden units: 32

Critic:



- Number of hidden units: 32

Note: Label "0.0001 128" means that the learning curve plots the result of the Actor-Critic algorithm with learning rate 0.0001 and minibatch size 128 for gradient descent. "buo" (baseline update once) means that baseline update was made through batch gradient descent instead of minibatch.

From the graphs we can see that in five out of six learning rate settings, using minibatch gradient descent achieves better result than simply performing batch gradient descent.

Figures 14 and 15 give a further comparison of different learning rates as well as mini-batch sizes. From the comparison, it seems that Actor-Critic algorithm with learning rate 0.001 and minibatch size 256 achieves the best performance among all tested settings.

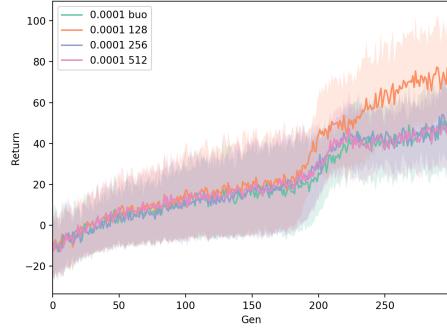


Figure 8

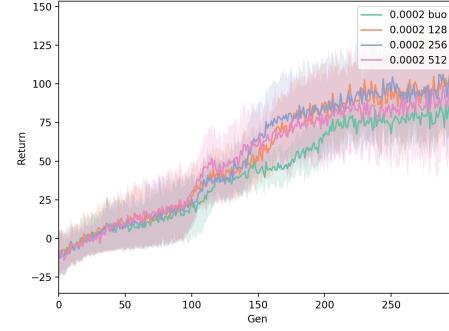


Figure 9

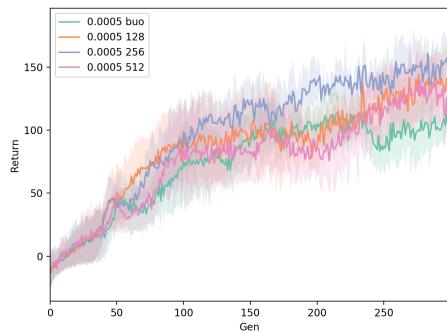


Figure 10

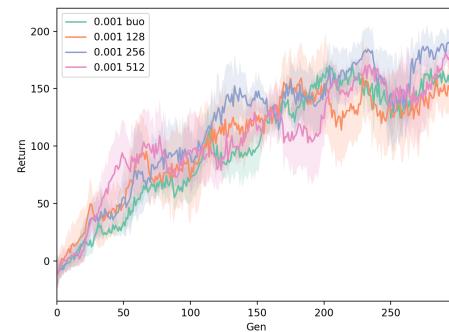


Figure 11

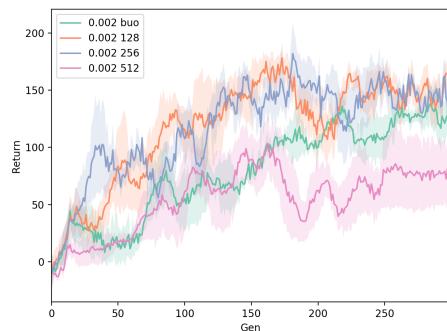


Figure 12

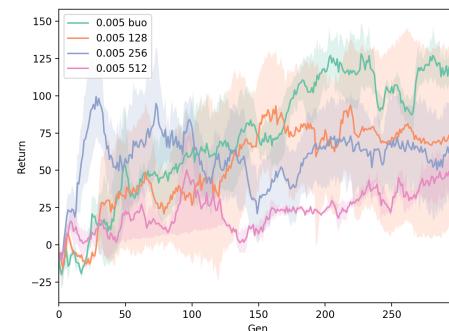


Figure 13

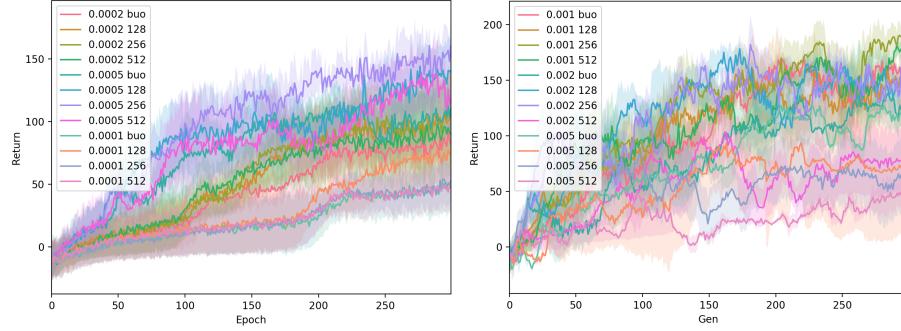


Figure 14

Figure 15

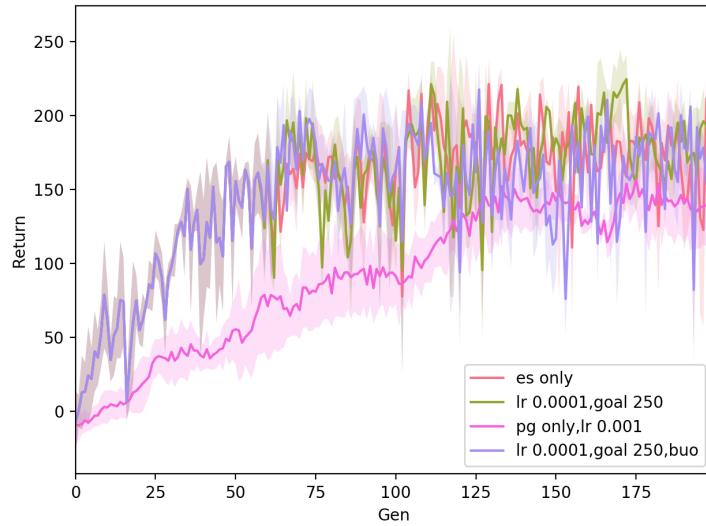


Figure 16

Figure 16 offers a comparison of the performances of Actor-Critic algorithm with learning rate 0.0001, vanilla Evolution Strategy, as well as Algorithm 2.

Note: "es only" is the learning curve for vanilla Evolution Strategy; "pg only, lr 0.001" is the learning curve for Actor-Critic algorithm with learning rate 0.001 and batch gradient descent (the vanilla Actor-Critic algorithm with the best performance); "lr 0.0001, goal 250" is the learning curve for Algorithm 2, with learning rate 0.0001, learning goal 250, and minibatch size 256; "lr 0.0001, goal 250, buo" is the learning curve for Algorithm 2, with learning rate 0.0001, learning goal 250, and batch gradient descent. The hyper-parameters for Evolution Strategy are the same as the ones specified in Section 1.2.

From the graph we see that vanilla Evolution Strategy still outperforms vanilla Actor-Critic algorithm. The performance of Algorithm 2 matches that of vanilla ES, and is

even more stable than that of vanilla ES once the threshold in this algorithm, 80% learning goal, is achieved. Comparing the performances of Algorithm 2 with and without minibatch gradient descent, it seems that minibatch gradient descent helps to improve the stability of performance.

3.1.2 Different learning rates for Actor and Critic

To further improve the performance of Algorithm 2, we experimented with using different learning rates for the Actor and Critic networks separately. Experiment results are shown in figures 17-22. Neural network architecture and hyper-parameters in these experiments are the same as in Section 2.1.1.

Note: "Alr 0.0001, Clr 0.0001" means that the learning curve plots the result of Algorithm 2, with Actor network learning rate 0.0001 and Critic network learning rate 0.0001. Figures 17-19 each compares results with one fixed Actor network learning rate and six different Critic network learning rates, while figures 20-22 each compares results with one fixed Critic network learning rates and six different Actor network learning rates

However, among the various different combinations of learning rates we have tested, none of them seems to achieve significantly better result than a shared learning rate between the Actor and Critic networks. Based on these results, we concluded that using different learning rates for the Actor and Critic networks separately is not very helpful in improving the performance of Algorithm 2.

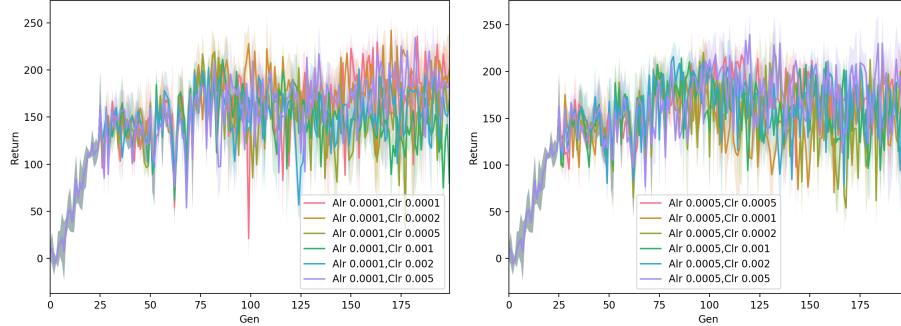


Figure 17

Figure 18

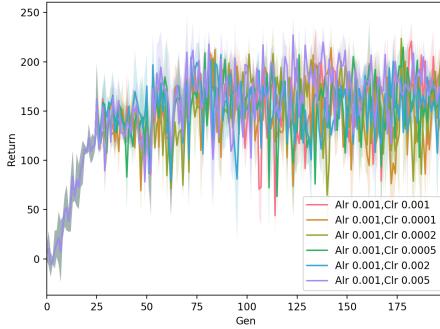


Figure 19

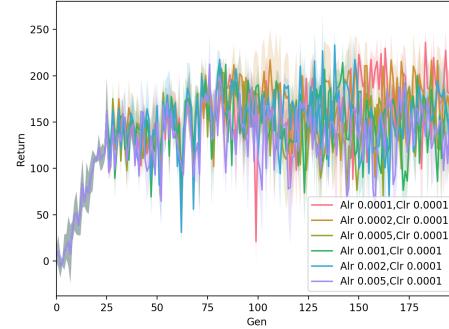


Figure 20

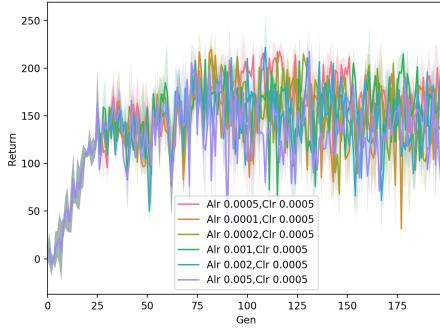


Figure 21

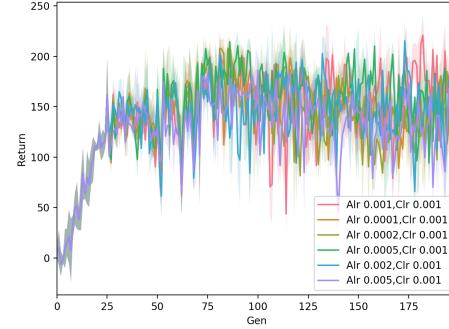


Figure 22

3.2 Combine Evolution Strategy with PPO

We also implemented Algorithm 2 with PPO algorithm, one of the most powerful Policy Gradient methods known now, in order to improve the performance of Algorithm 2. Since PPO algorithm has the reputation for being sample efficient, this PPO version of Algorithm 2 also has the potential to outperform vanilla Evolution Strategy in terms of sample efficiency. Among the environments we tested, the sample efficiency (amount of data generated when running a fixed number of generations/iterations) of PPO algorithm is always better than that of Evolution Strategy. Thus we expected Algorithm 2 to be more sample efficient than vanilla Evolution Strategy as well.

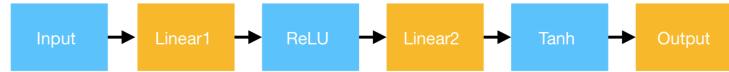
The hyper-parameter setting of PPO algorithm as well as neural network architecture is shared in all the following experiments, and is listed below:

- episode length: 100

- minibatch size per update: 64
- number of epoch per iteration: 10
- learning rate: 3e-4
- weight decay rate: 1e-3
- lambda value in generalized advantage estimator: 0.95
- Neural network architecture:
Policy network:



Baseline network:



3.2.1 Experiment results on InvertedPendulum environment

Below are the additional hyper-parameters we used on the InvertedPendulum-v2 environment:

- α : 5
- σ : 1
- Number of hidden units: 32
- Number of generations: 50

We experimented with different population sizes for Evolution Strategy as well as different number of episodes per PPO iteration. Results are shown in figures 23 and 24. Note: "IP_es_pop100_505000" means that the learning curve plots the result of vanilla ES with population size 100 on the InvertedPendulum-v2 environment, and the total time steps of this experiment is 505000. "IP_pg_nep21_105000" means that the learning curve plots the result of PPO algorithm with 21 episodes generated per iteration, and the total time steps is 105000. All results were averaged over three random seeds. The graphs imply that enlarging population size in Evolution Strategy and increasing the number of episodes generated per PPO iteration both have a positive effect on the performances of the algorithms. Comparing the two graphs, it's clear that the learning curves of vanilla ES grow much faster than those of PPO, especially at the beginning of the experiments. However, the curves of PPO are much more smooth than those of ES, implying a more stable policy performance improvement and a stronger tendency to converge. What's more, the amount of data generated in each experiment using PPO

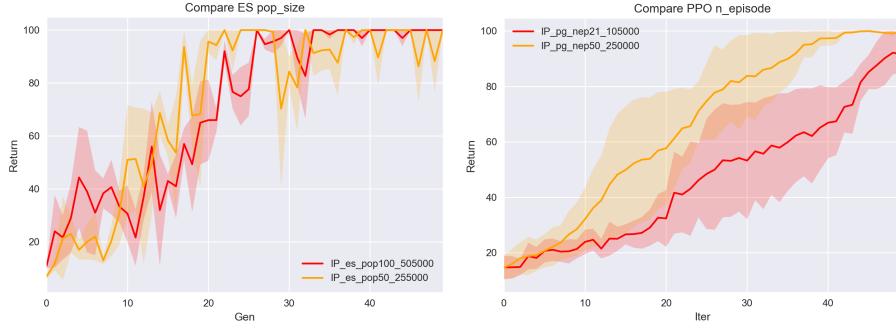


Figure 23

Figure 24

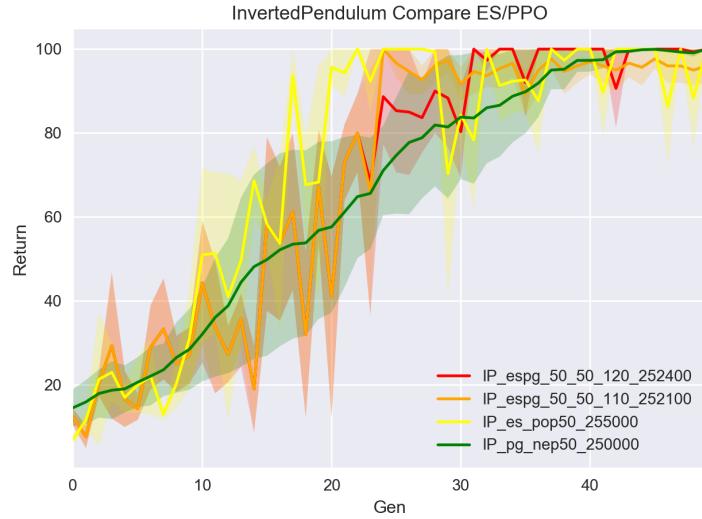


Figure 25

algorithm is generally less than that of vanilla ES.

Note: "IP_espg_50_50_120_252400" means that the learning curve plots the result of Algorithm 2, with population size 50 per ES generation and 50 episodes per PPO iteration, learning goal 120, and total time steps 252400 per experiment.

Figure 25 compares the performances of vanilla ES, PPO, and Algorithm 2 on the InvertedPendulum-v2 environment. The results were averaged over three random seeds. In our experiments, the maximum episode length was set to be 100. Since the maximum reward in this environment with episode length 1000 is 1000, we set the learning goal to be 120 and 110 in our experiments. From the results it's hard to tell whether Algorithm 2 performs better than vanilla ES, which is probably due to the fact that it takes PPO more iterations to achieve result that is as good as vanilla ES, and thus the

transition from ES to PPO is not so smooth. We also note that at the end of 50 generations, the final result of the experiment with learning goal 110 is still slightly lower than the result of vanilla ES as well as Algorithm 2 with learning goal 120. We then speculate that if the threshold is set too low and we switch to PPO too soon, the current policy and the real optimum will still be far and the improvement afterwards will be slow.

3.2.2 Experiment results on HalfCheetah environment

Below are the additional hyper-parameters we used on the HalfCheetah-v2 environment:

- α : 0.1
- σ : 0.05
- Number of hidden units: 64
- Number of generations: 150

Comparisons of different population sizes in Evolution Strategy and different number of episodes per PPO iteration are shown in figures 26 and 27.

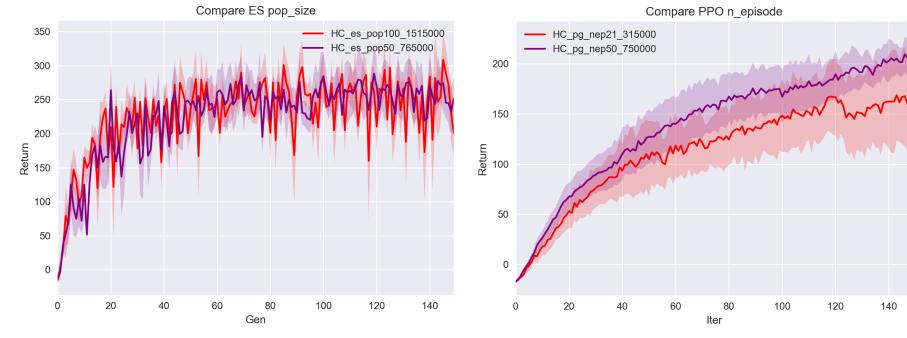


Figure 26

Figure 27

Figures 28 and 29 compare the performances of vanilla ES, PPO, and Algorithm 2 on the HalfCheetah-v2 environment. The results were averaged over three random seeds. In our experiments, the maximum episode length was set to be 100. Since the current benchmark reward for this environment with maximum episode length 1000 is around 3000, we conducted experiments with learning goals 375 and 350. In both graphs, learning curves for Algorithm 2 (the red, yellow, and orange ones, plotted with different pop_size/n_episode ratios) generally grow as fast as that for vanilla ES, but are less volatile. From the curves of different pop_size/n_episode ratios, we can see that larger population size tends to generate less volatile result. We also note that the total times steps of Evolution Strategy doubles that of PPO, and that integrating PPO into Evolution Strategy indeed helps to improve sample efficiency of vanilla ES.

Note: "HC_esp50_50_375_761733" means that the learning curve plots the result of Algorithm 2, with population size 50 per ES generation and 50 episodes per PPO iteration, learning goal 375, and total time steps 761733 per experiment.

Figure 30 compares the results of vanilla ES and Algorithm 2 with learning goals 350 and 375. Learning goal doesn't seem to affect the performance of the algorithm much in this case. However, figure 30 shows that the total time steps in the experiment with learning goal 350 is less than that with learning goal 375 by 15%. We then speculate that by reasonably lowering the learning goal and thus increasing PPO iterations, we can see further drop in the amount of data used.

3.2.3 Experiment results on Hopper environment

The hyper-parameters on Hopper-v2 environment are the same as the ones on HalfCheetah-v2 environment.

Note: "HP_esp50_50_250966" means that the learning curve plots the result of Algorithm 2, with population size 50 per ES generation and 50 episodes per PPO iteration, and total time steps 250966 per experiment.

Figure 31 compares different population sizes in Evolution Strategy and figure 32 compares different numbers of episodes per PPO iteration. The two figures show that both vanilla ES and PPO can get good results within few generations/iterations on this environment.

Figure 33 compares the performances of vanilla ES, PPO, and Algorithm 2 on the Hopper-v2 environment. The results were averaged over three random seeds. In our experiments, the maximum episode length was set to be 100. Since the results of both vanilla ES and PPO became stable around 220 after 50 generations/iterations, we ran our new algorithm for 50 generations and set the learning goal to be 265. Similar to the HalfCheetah environment, vanilla ES learns faster than PPO at the beginning of the experiments. But the stabilizing effect of PPO in Algorithm 2 is clearer on this environment. From our understanding, this improvement in stabilization is resulted from the fact that PPO is able to achieve performance as strong as that of vanilla ES within a similar number of iterations on this environment.

4 Possible Future Work

- Algorithm 1 should be further tested by combining ES with PPO
- Both Algorithm 1 and Algorithm 2 need to be tested on more environments, in both MuJoCo and Atari
- Experiments should be re-run with the maximum episode length of each environment
- Sample efficiency of different algorithms need to be compared in a more rigorous fashion

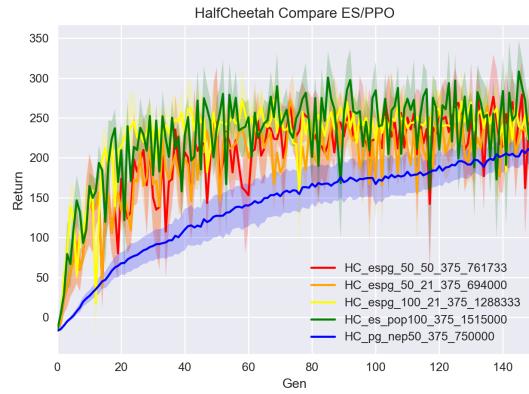


Figure 28

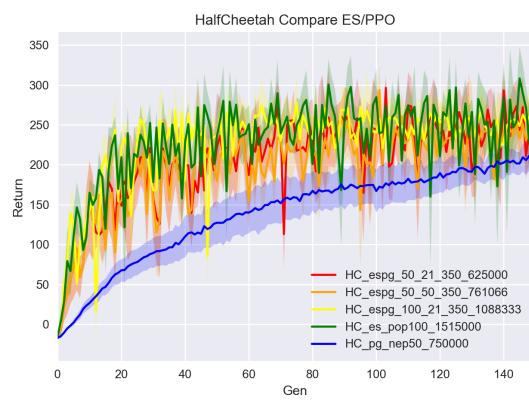


Figure 29

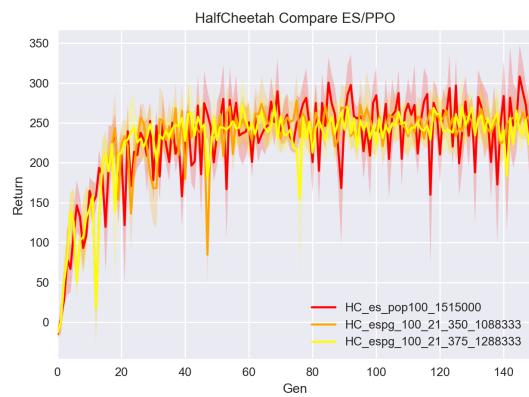


Figure 30

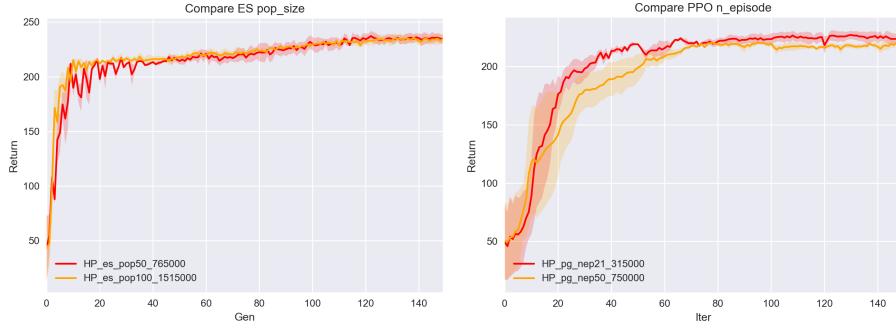


Figure 31

Figure 32

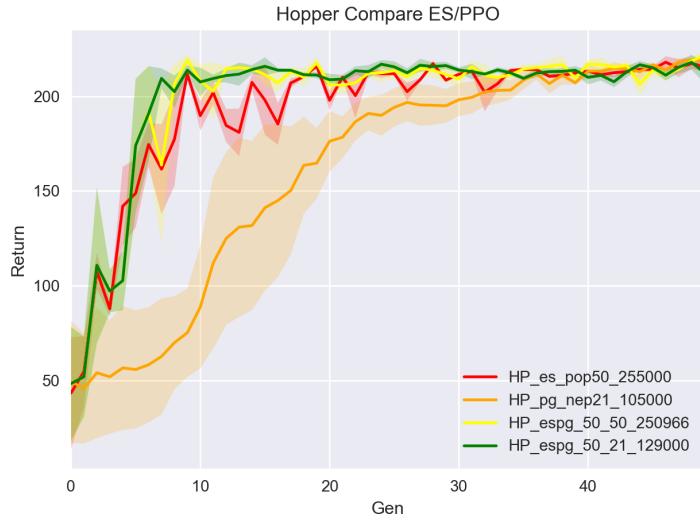


Figure 33

References

- [1] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pages 1057–1063, 2000.
- [2] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International Conference on Machine Learning*, pages 1889–1897, 2015.
- [3] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- [4] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [5] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- [6] 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.