

403 A Derivation

404 The AMT and ES update rules are given as

$$\begin{aligned} \sigma_{(t)} &\leftarrow \sigma_{(t-1)} + \frac{\alpha_{es}}{n\sigma_{(t-1)}} \text{SmoothL1}(R_{max,(t-1)}, R_{avg,(t-1)}) \\ \theta_{(t+1)} &\leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\sigma_{(t)}} \sum_{i=1}^N R_i \epsilon_i \end{aligned}$$

406 Using the expression for $\sigma_{(t)}$ in the ES update yields the following

$$\begin{aligned} \theta_{(t+1)} &\leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n(\sigma_{(t-1)} + \frac{\alpha_{es}}{n\sigma_{(t-1)}} \text{SmoothL1}(R_{max,(t-1)}, R_{avg,(t-1)}))} \sum_{i=1}^n R_i \epsilon_i \\ &= \theta_{(t+1)} \leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\sigma_{(t-1)}(1 + \frac{\alpha_{es}}{n\sigma_{(t-1)}^2} \text{SmoothL1}(R_{max,(t-1)}, R_{avg,(t-1)}))} \sum_{i=1}^n R_i \epsilon_i \\ &= \theta_{(t+1)} \leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\sigma_{(t-1)}\Lambda_{(t-1)}} \sum_{i=1}^n R_i \epsilon_i \end{aligned}$$

408 where $\Lambda_{(t-1)} = 1 + \frac{\alpha_{es}}{n\sigma_{(t-1)}^2} \text{SmoothL1}(R_{max,(t-1)}, R_{avg,(t-1)})$. Expanding $\sigma_{(t-1)}$ using the AMT update

409 rule gives us

$$\begin{aligned} \theta_{(t+1)} &\leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\Lambda_{(t-1)}(\sigma_{(t-2)} + \frac{\alpha_{es}}{n\sigma_{(t-2)}} \text{SmoothL1}(R_{max,(t-2)}, R_{avg,(t-2)}))} \sum_{i=1}^n R_i \epsilon_i \\ &= \theta_{(t+1)} \leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\Lambda_{(t-1)}\sigma_{(t-2)}(1 + \frac{\alpha_{es}}{n\sigma_{(t-2)}^2} \text{SmoothL1}(R_{max,(t-2)}, R_{avg,(t-2)}))} \sum_{i=1}^n R_i \epsilon_i \\ &= \theta_{(t+1)} \leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\Lambda_{(t-1)}\sigma_{(t-2)}\Lambda_{(t-2)}} \sum_{i=1}^n R_i \epsilon_i \end{aligned}$$

411 where $\Lambda_{(t-2)} = 1 + \frac{\alpha_{es}}{n\sigma_{(t-2)}^2} \text{SmoothL1}(R_{max,(t-2)}, R_{avg,(t-2)})$. Expanding this recursively gives us the

412 following

$$\begin{aligned} \theta_{(t+1)} &\leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\sigma_{(1)}\Lambda_{(t-1)}\Lambda_{(t-2)}\dots\Lambda_{(1)}} \sum_{i=1}^n R_i \epsilon_i \\ &= \theta_{(t+1)} \leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\sigma_{(1)} \prod_{t'=1}^{t-1} \Lambda_{(t')}} \sum_{i=1}^n R_i \epsilon_i \\ &= \theta_{(t+1)} \leftarrow \theta_{(t)} + \frac{\alpha_{es}}{n\sigma_{(1)}\hat{\Lambda}} \sum_{i=1}^n R_i \epsilon_i \end{aligned}$$

414 Hence, yielding the AMT update in the form of initial mutation rate $\sigma_{(1)}$.

415 B Policy Improvement

416 We will now show that AMT improves the policy of winners in the population. Let us consider two successive
417 gradient intervals g indexed by (l) and $(l-1)$. Let $p_{(l)}$ and $p_{(l-1)}$ be the probabilities of convergence to the
418 optimal policy $\pi_{\theta_{es}}^*(a_t|s_t)$ in the weight space at (l) and $(l-1)$ respectively.

420 We start by evaluating the mutation rates at (l) and $(l-1)$ which are given as $\sigma_{(l)} > \sigma_{(l-1)}$. We can now
421 evaluate the probabilities of convergence to $\pi_{\theta_{es}}^*(a_t|s_t)$ as

$$p_{(l)} \geq p_{(l-1)}$$

422 Using this fact, we can evaluate the winners (indexed by q) in the sorted reward population F .

$$\begin{aligned}
\sum_{q=1}^w p_{(l)}^{(q)} F_{(l)}^{(q)} &\geq \sum_{q=1}^w p_{(l-1)}^{(q)} F_{(l-1)}^{(q)} \\
&= \mathbf{E}_{F_{(l)}^{(q)} \sim F_{(l)}} [F_{(l)}^{(q)}] \geq \mathbf{E}_{F_{(l-1)}^{(q)} \sim F_{(l-1)}} [F_{(l-1)}^{(q)}] \\
&= \mathbf{E}[W_{(l)}] \geq \mathbf{E}[W_{(l-1)}]
\end{aligned}$$

423 Here, $p_{(l)}^{(q)}$ is the probability of convergence of actor q (having observed reward $F_{(l)}^{(q)}$) to its optimal policy
424 $\pi_{\theta_{es}}^{(q),*}(a_t|s_t)$ at interval (l) . $W_{(l)}$ represents the set of winners at (l) . The mathematical expression obtained
425 represents that the set of winners $W_{(l)}$ formed at the next gradient interval (l) is at least as good as the previous
426 set of winners $W_{(l-1)}$, i.e.- $\pi_{\theta_{es},(l)}^{(q)}(a_t|s_t) \geq \pi_{\theta_{es},(l-1)}^{(q)}(a_t|s_t)$. This guarantees policy improvement among
427 winners of the population.

428 C Additional Results

429 C.1 Performance

430 We evaluate the performance and sample-efficiency of ESAC on 21 MuJoCo and DeepMind Control Suite
431 [38]. Figure Figure 7 presents learning behavior of ESAC in comparison to SAC, TD3, PPO and ES on all
432 21 tasks. Training setup for all agents was kept same with different values of hyperparameters (presented in
433 subsection D.2). ESAC demonstrates improved returns on 16 out of 21 tasks.

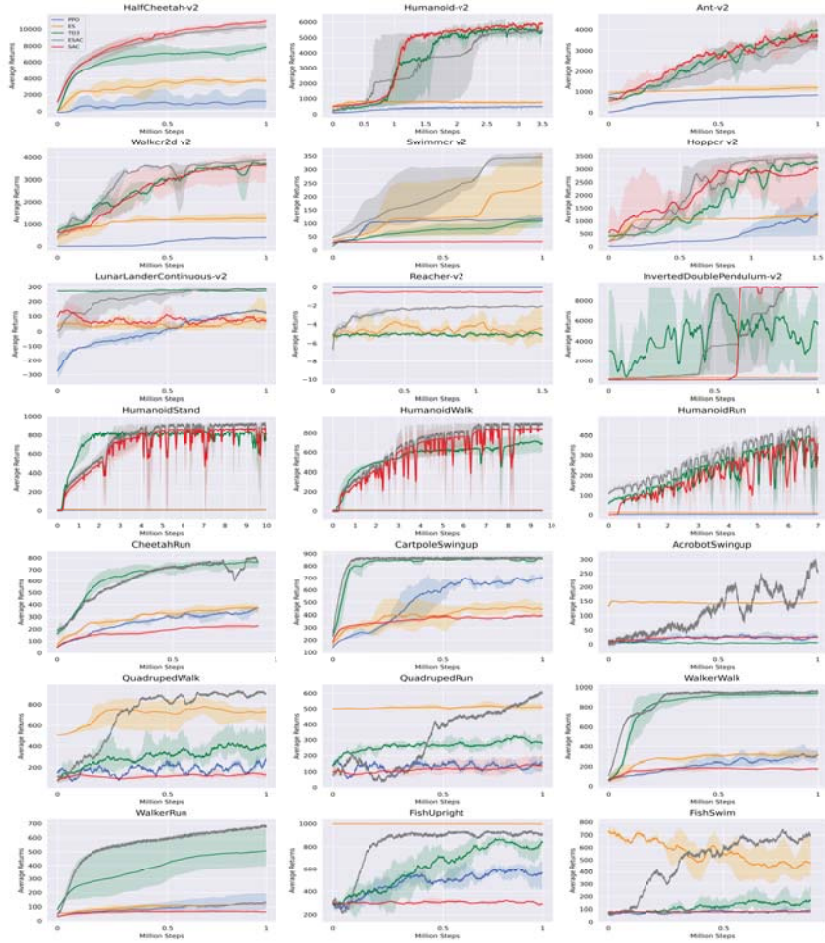


Figure 7: Average Returns on MuJoCo and DeepMind Control Suite tasks. ESAC’s demonstrates state-of-the-art performance on out of 21 tasks.

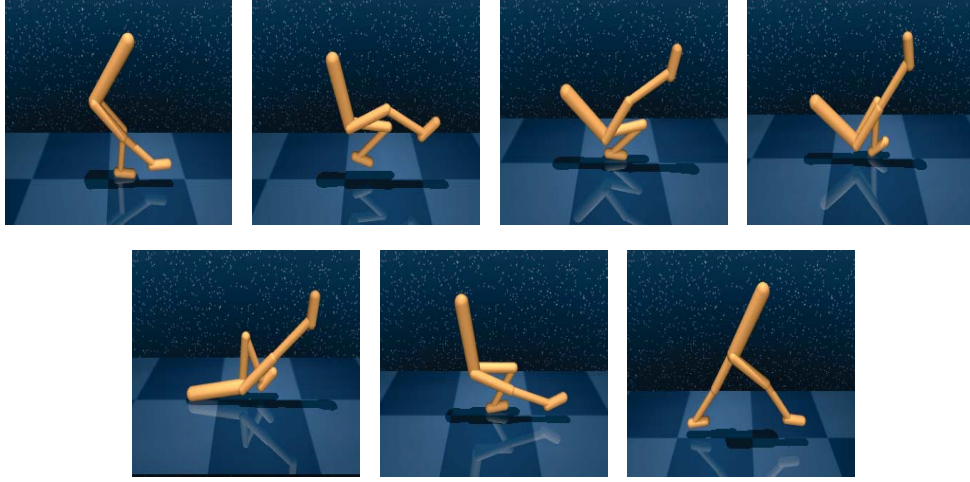


Figure 8: Robust behavior of ESAC observed on the WalkerWalk task. The ESAC policy prevents the walking robot from falling down when the robot loses its balance while walking. Although ESAC falls short of optimal performance as a result of sparse rewards, it exhibits robust policies on complex tasks due to successive evolutions.

434 C.2 Scalability

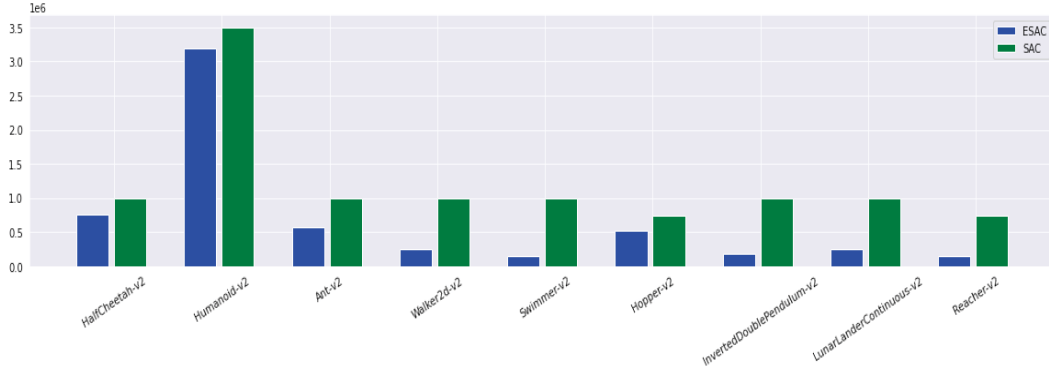


Figure 9: Complete results on the number of backprop updates for MuJoCo control tasks and LunarLanderContinuous environment from OpenAI’s Gym suite. ESAC exponentially anneals gradient-based SAC updates and leverages winner selection and crossovers for computational efficiency and performance improvement.

435 D Hyperparameters

436 D.1 MuJoCo

437 Hyperparameter values for our experiments are adjusted on the basis of complexity and reward functions
 438 of tasks. In the case of MuJoCo control tasks, training-based hyperparameters are kept mostly the same
 439 with the number of SAC episodes varying as per the complexity of task. Tuning was carried out on
 440 HalfCheetah-v2, Ant-v2, Hopper-v2 and Walker2d-v2 tasks. Out of these, Ant-v2 presented high variance
 441 indicating the requirement of a lower learning rate. Number of SAC updates in SAC and ESAC implemen-
 442 tations were kept 1 for a fair comparison with TD3. All tasks have a common discount factor $\gamma = 0.99$,
 443 SAC learning rate $\alpha = 3 \times 10^{-4}$, population size $n = 50$, mutation rate $\sigma = 5 \times 10^{-3}$ and winner frac-
 444 tion $e = 0.4$. ES learning rate α_{es} was kept fixed at 5×10^{-3} for all tasks except Ant-v2 having $\alpha_{es} = 1 \times 10^{-4}$.
 445

446 The only variable hyperparameter in our experiments is number of SAC episodes executed by the SAC agent.
 447 Although the ESAC population in general is robust to its hyperparameters, the SAC agent is sensitive to the
 448 number of episodes. During the tuning process, the number of episodes were kept constant to a value of 10 for
 449 all the tasks. However, this led to inconsistent results on some of the environments when compared to SAC
 450 baseline. As a result, tuning was carried out around this value to obtain optimal results corresponding to each
 451 task. The SAC agent executed a total of 10 episodes for each gradient interval g for Ant-v2, Walker2d-v2,
 452 Hopper-v2 and Humanoid-v2 tasks; 5 episodes for HalfCheetah-v2, LunarLanderContinuous-v2, Reacher-v2
 453 and InvertedPendulum-v2 tasks; and 1 episode for Swimmer-v2 task.

454 **D.2 DeepMind Control Suite**

455 The DeepMind control suite presents a range of tasks with sparse rewards and varying complexity for the same
 456 domain. Hyperparameter values for these tasks are different from that of MuJoCo control tasks. Tuning was
 457 carried on the Cheetah, Quadruped and Walker tasks in order to obtain sample-efficient convergence. In the case
 458 of SAC, granularity of hyperparameter search was refined in order to observe consistent behavior. However,
 459 different values of temperature parameter produced varying performance. As in the MuJoCo case, we kept the
 460 number of updates fixed to 1 in order to yield a fair comparison with TD3.

461
 462 All tasks have a common discount factor $\gamma = 0.99$, SAC learning rate $\alpha = 3 \times 10^{-4}$, population size $n = 50$,
 463 mutation rate $\sigma = 1 \times 10^{-2}$, winner fraction $e = 0.4$ and ES learning rate 1×10^{-2} . As in the case of MuJoCo
 464 tasks, the SAC is found to be sensitive to the number of episodes during the gradient interval g . These were kept
 465 constant at 5 and then tuned around this value for optimal performance. Final values of the SAC episodes were 5
 466 for CartpoleSwingup, WalkerWalk and WalkerRun tasks and 1 for the remaining tasks.