

---

# The Cross-environment Hyperparameter Setting Benchmark for Reinforcement Learning

---

Andrew Patterson, Samuel Neumann, Raksha Kumaraswamy, Martha White, Adam White  
Department of Computing Science, University of Alberta  
{ap3,sfneuman,kumarasw,whitem,amw8}@ualberta.ca

## Abstract

1 This paper introduces a new benchmark, the Cross-environment Hyperparameter  
2 Setting Benchmark, that allows comparison of RL algorithms across environments  
3 using only a single hyperparameter setting, encouraging algorithmic development  
4 which is insensitive to hyperparameters. We demonstrate that the benchmark is  
5 robust to statistical noise and obtains qualitatively similar results across repeated  
6 applications, even when using a small number of samples. This robustness makes  
7 the benchmark computationally cheap to apply, allowing statistically sound insights  
8 at low cost. We provide two example instantiations of the CHS, on a set of six  
9 small control environments (SC-CHS) and on the entire DM Control suite of 28  
10 environments (DMC-CHS). Finally, to demonstrate the applicability of the CHS to  
11 modern RL algorithms on challenging environments, we provide a novel empirical  
12 study of an open question in the continuous control literature. We show, with  
13 high confidence, that there is no meaningful difference in performance between  
14 Ornstein-Uhlenbeck noise and uncorrelated Gaussian noise for exploration with  
15 the DDPG algorithm on the DMC-CHS.

## 16 1 Introduction

17 One of the major benefits of the Atari suite is the focus on more general reinforcement learning agents.  
18 Numerous agents have been shown to exhibit learning across many games with a single architecture  
19 and a single set of hyperparameters. To a lesser extent, OpenAI Gym (Brockman et al., 2016) and DM  
20 control suite (Tassa et al., 2018) are used in the same way—though at times not all environments are  
21 used, raising the possibility of cherry-picking. As the ambitions of the community have grown, Atari  
22 and OpenAI Gym tasks have been combined into larger problem suites, with subsets of environments  
23 chosen to test algorithms. In many ways we are back to where we started with Cartpole, Mountain  
24 Car and the like: where environment-specific hyperparameter tuning and problem subselection is  
25 prominent. Instead of proposing a new and bigger challenge suite, we explore a challenging new  
26 benchmark and empirical methodology for comparing agents across a given set of environments,  
27 complementing the existing empirical toolkit for investigating the scalability of deep RL algorithms.

28 In order to make progress towards impactful applications of reinforcement learning and the broader  
29 goals of AGI, we need benchmarks that clearly highlight the generality and stability of learning  
30 algorithms. Empirical work in Atari, Mujoco, and simulated 3D worlds typically use networks with  
31 millions of parameters, dozens of GPUs, and up to billions of samples (Beattie et al., 2016; Espeholt  
32 et al., 2018). Many results are demonstrative, meaning that the primary interest is not the stability  
33 and sensitivity, nor what was required to achieve the result, rather that the result *could* be achieved.  
34 It is infeasible to combine these large scale experiments with hyperparameter studies and enough  
35 independent runs to support statistically significant comparisons. More evidence is emerging that  
36 such state of the art systems (1) rely on environment-specific design choices that are sensitive to  
37 minor changes to hyperparameters (Henderson et al., 2018; Engstrom et al., 2019), (2) are less data

38 efficient and stable compared with simple baselines (van Hasselt et al., 2019; Taïga et al., 2019),  
 39 and (3) cannot solve simple toy tasks without extensive re-engineering (Obando-Ceron and Castro,  
 40 2021; Patterson et al., 2021). It is abundantly clear that modern RL methods can be adapted to a  
 41 broader spectrum of challenging tasks—well beyond what was possible with linear methods and  
 42 expert feature design. However, we must now progress to phase two of empirical deep RL research:  
 43 focusing on generality and stability.

44 There is a growing movement to increase the standards of empirical work in RL. Noisy results,  
 45 inconsistent evaluation practices, and divergent code bases have fueled calls for more open-sourcing  
 46 of agent architecture code, experiment checklists, and doing more than three independent evaluations  
 47 in our experiments (Henderson et al., 2018; Pineau et al., 2020). Digging deeper, recent work has  
 48 highlighted our poor usage of basic statistics, including confidence intervals and hypothesis tests  
 49 (Colas et al., 2018). Long before the advent of deep networks, researchers called out the environment  
 50 overfitting that is rampant in RL and proposed sampling from parameterized variants of classic  
 51 control domains to emphasize general methods (Whiteson et al., 2009). Finally, and most related to  
 52 our work, Jordan et al. (2020) proposed a methodology to better characterize the performance of an  
 53 algorithm across environments, evaluated with randomly sampled hyperparameters. We build on this  
 54 direction, but focus on a simpler and more computationally frugal evaluation that examines the single  
 55 best hyperparameter setting across environments, rather than a randomly sampled one, and allows for  
 56 a smaller number of runs per environment.

57 Table 1: Chance of incorrect claims

	3 runs	10	30	100
58 Acrobot	47%	31%	22%	1%
59 Cartpole	7%	0%	0%	0%
60 CliffWorld	54%	19%	14%	0%
61 LunarLander	16%	7%	1%	0%
62 MountainCar	22%	9%	7%	0%
63 PuddleWorld	18%	16%	8%	0%

64 One reason we focus on computational efficiency is that computational limitations seems to be the pri-  
 65 mary culprit for misleading or incorrect claims in RL experiments. Experiments with many runs, many hy-  
 66 perparameters, and many environments can be computationally prohibitive. The typical trade-off is to  
 67 use a smaller number of runs. Such a choice, however, can lead to incorrect conclusions. Table 1 shows  
 68 the empirical probability of incorrectly ordering four reasonable RL algorithms across several domains often considered too small to draw meaningful con-  
 69 clusions. We ran each of the four algorithms 250 times on every domain and for every hyperparameter  
 70 setting in an extensive sweep to get a high confidence approximation of the correct ordering between  
 71 algorithms. We then used bootstrap sampling to simulate 10k papers—each using a small number of  
 72 random seeds—and counted the frequency that incorrect algorithm orderings were reported. Even  
 73 with 30 runs in these small domains, incorrect rankings were **not** uncommon. Further details are  
 74 described in Section 5.

75 Another critical issue for algorithm evaluation is the difficulty in hyperparameters selection. Modern  
 76 RL algorithms require tuning an increasing number of hyperparameters, greatly impacting the  
 77 outcome of an experimental trial. As more hyperparameters are introduced, the computational  
 78 burden of tuning grows exponentially. To combat this, several strategies have emerged in the literature  
 79 including relying on default hyperparameter values (Schaul et al., 2016; Wang et al., 2016; Van Hasselt  
 80 et al., 2016), tuning hyperparameters on a subset of domains (Bellemare et al., 2013), or eroding  
 81 standards of sufficient statistical power for publication (Henderson et al., 2018; Colas et al., 2018).

82 Our new benchmark is designed to (1) standardize the selection of hyperparameters, (2) evaluate  
 83 stability over runs, (3) be computationally cheap to run, and (4) be easy to use. We propose the Cross-  
 84 environment Hyperparameter Setting Benchmark (CHS). The basic idea is simple: an algorithm is  
 85 evaluated on a set of environments, using the best hyperparameter setting across those environments,  
 86 rather than per-environment. Though conceptually simple, this methodology is not widely used.  
 87 We first address some of the nuances in the CHS, namely how to standardize performance across  
 88 environments to allow for aggregation, how to allow for robust measures of performance, and finally  
 89 how to reduce computation to make it more feasible to use the CHS. We evaluate the effectiveness of  
 90 the CHS itself by examining the stability of the conclusions from the CHS under different numbers  
 91 of runs. We then demonstrate that the CHS can result in different conclusions about algorithms  
 92 compared to the conventional *per-environment tuning* approach and the more recent approach of  
 using a subset of environments for tuning. We conclude with a larger demonstration of the CHS on  
 DM Control Suite.

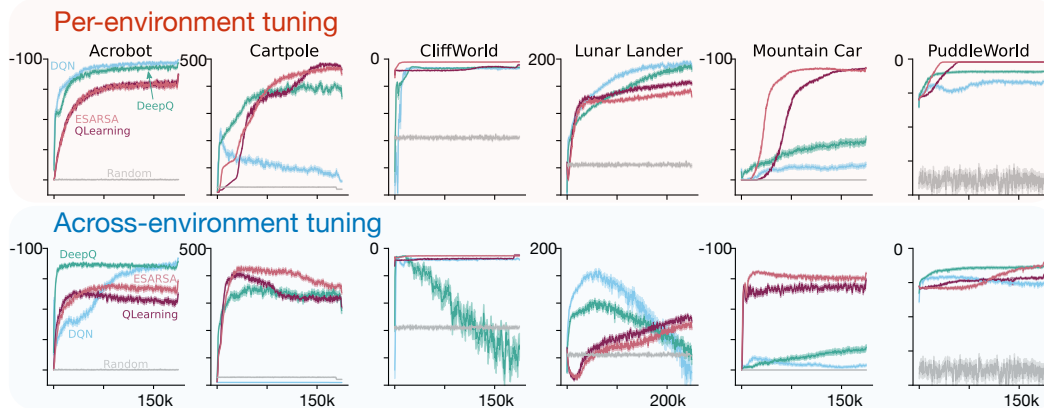
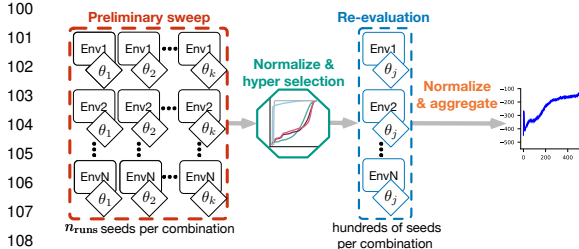


Figure 1: An example experiment comparing four algorithms across six different environments. Each learning curve shows the mean and standard error of 250 independent runs for each algorithm and environment. Hyperparameters are selected using three runs of every algorithm, environment, and hyperparameter setting. **Top** shows the learning curves when the best hyperparameters are chosen for each environment individually. **Bottom** shows the learning curves when hyperparameters are chosen according to our benchmark, the CHS.

## 93 2 Contrasting Across-Environment versus Per-Environment Tuning

94 In this section, we introduce the basic procedure for the CHS and provide an experiment showing  
 95 how it can significantly change empirical outcomes compared to the conventional per-environment  
 96 tuning approach. We provide specific details for each step later and here focus on outlining the basic  
 97 idea and its utility.

98 The CHS consists of the following four steps summarized in the inset figure below. We assume we  
 99 are given a set of environments and a set of hyperparameters for the algorithm we are evaluating.



**Step 1 (Preliminary Sweep)** Run the algorithm for all hyperparameters and all environments, for  $n_{\text{runs}}$  runs (i.e.,  $n_{\text{runs}} < 30$ ) and record the performance of every combination. The performance could be online average return per step.

**Step 2 (Normalization)** Normalize the scores across environments to be in  $[0, 1]$ . We use CDF normalization, which is described in Section 4.

**Step 3 (Hyperparameter Selection)** Select the hyperparameter setting with the highest score averaged across environments.

109 **Step 4 (Re-evaluation)** With the single best hyperparameter setting, use many more runs in each  
 110 environment (e.g. 100) to produce a more accurate estimate of performance.  
 111

112 The last step is more lightweight than it appears since only a single hyperparameter setting is used  
 113 for all environments. Executing 100 or more runs for every hyperparameter setting would likely be  
 114 prohibitive. The trick is to use a small  $n_{\text{runs}}$  in the Preliminary Sweep, saving compute, and a larger  
 115 number of runs in the Re-evaluation step. This contrasts the conventional *per-environment tuning*  
 116 approach of choosing hyperparameter settings which maximize performance on each environment  
 117 individually, which can be more sensitive when  $n_{\text{runs}}$  is small.

118 We now show an experiment comparing the CHS and this conventional approach in Figure 1. The  
 119 per-environment tuning approach highlights the ideal behavior of an algorithm per environment,  
 120 whereas the CHS highlights the (in)sensitivity of an algorithm across environments. Experimental  
 121 details can be found in Section 5. The environments are relatively simple (most coming from the  
 122 classic control suite of OpenAI Gym (Brockman et al., 2016)) but difficult enough for our purposes:  
 123 no one algorithm could reach near optimal performance in all environments.

124 The CHS does not rank the algorithms differently than with per-environment tuning, but CHS  
 125 does alert us to potential catastrophic failure of some algorithms. The neural network DeepQ  
 126 agent performs terribly in Cliffworld and Lunar Lander under CHS, but appears reliable under  
 127 the per-environment approach. What is going on? Forced to select only one hyperparameter

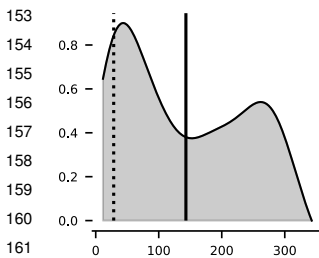
128 across environments, the best outcome is to sacrifice performance in Cliffworld and Lunar Lander—  
129 achieving worse performance than a uniform random policy.

### 130 3 Performance Distributions

131 In this section, we describe the distribution and random variables underlying an RL experiment. This  
132 formalism allows us to reason about the summary statistics we consider for the CHS in the next  
133 section. We also visualize these distributions to provide intuition on the properties of the summary  
134 statistics of these distributions and the implications for the single performance numbers used in RL.

135 In an RL experiment, we seek to describe the performance distribution of an algorithm for each  
136 hyperparameter setting  $\theta \in \Theta$ , denoted as  $\mathbb{P}(G, E | \theta)$  where  $G$  is a random variable indicating  
137 the performance of an algorithm on a given environment,  $E \in \mathcal{E}$ . Most commonly, we report  
138 an estimate of the average performance conditioned on environment and hyperparameter setting,  
139  $g(E, \theta) \cong \mathbb{E}[G | E, \theta]$  using a sample average and some measure of uncertainty about how accurately  
140  $g(E, \theta)$  approximates  $\mathbb{E}[G | E, \theta]$ .

141 The environment can be seen as a random variable for many RL experiments. The most common  
142 case is to specify a set of MDPs that the authors believe represent the important applications of their  
143 new algorithm. If results are uniformly aggregated across these environments, then this corresponds  
144 to assuming a uniform distribution over this set of environments. Other times, random subsets of  
145 environments from environment suites are chosen; the performance estimate on this subset provides  
146 an estimate of performance across the entire suite. The idea of evaluating algorithms over a random  
147 sample of MDPs has been studied explicitly previously. For example, the parameters determining  
148 the physics of classical control domains were randomized and sampled to avoid domain overfitting  
149 (Whiteson et al., 2009), and randomly generated MDPs (Archibald et al., 1995) have been used to  
150 evaluate new algorithmic ideas (Seijen and Sutton, 2014; Mahmood et al., 2014; White and White,  
151 2016). If we subselect after running the algorithms, then we bias the distribution over environments  
152 towards those with higher performance.



153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163 Figure 2: Performance  
164 distribution  $\mathbb{P}(G | E, \theta)$   
165 on Cartpole with hyper-  
166 parameter  $\text{stepsize}=2^{-9}$ .

Let us look at an example of these performance distributions to gain some  
intuition for estimating statistics like the expected performance. Con-  
sider the action-value nonlinear control method DQN, using the Adam  
optimizer (Mnih et al., 2013; Kingma and Ba, 2015), on Cartpole (Barto  
et al., 1983). We fix the hyperparameter setting  $\theta$  to the default values  
from Raffin et al. (2019). For this fixed environment, all randomness  
is due to sampling algorithm performance on this environment, namely  
sampling  $G$  according to  $\mathbb{P}(G | E, \theta)$ . The performance,  $G$ , is the average  
episodic return over all episodes completed during 100k learning steps.  
This environment is considered solved for  $G > 400$ . We repeat this pro-  
cedure for 250 independent trials to estimate the distribution  $\mathbb{P}(G | E, \theta)$ ,  
shown in Figure 2, with x-axis possible outcomes of  $G$  and y-axis the  
probability density. The vertical solid line denotes mean performance,  
and the vertical dotted line denotes mean performance of a random policy.

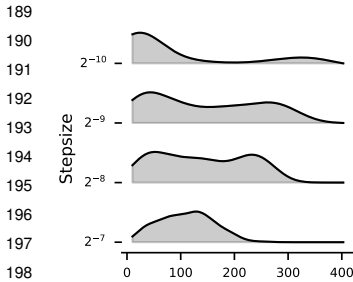
167 Figure 2 is a typical example of the performance of an RL algorithm over multiple independent  
168 trials. In this case, DQN is more likely to fail than to learn a policy which solves this relatively  
169 simple environment. It is common practice to run an RL algorithm for some number of random  
170 seeds—effectively drawing samples of performance from this distribution—then reporting the mean  
171 over those samples (solid vertical line).

172 There are two implications from observing this bimodal performance distribution. First, using the  
173 expected value of this distribution as the summary statistic does not aptly demonstrate that the poor  
174 performance of DQN on Cartpole is due to occasional catastrophic failure—performing worse than  
175 or equivalent to a random policy. Instead, mean performance might lead us to wrongly conclude that  
176 DQN on Cartpole usually finds a sub-optimal, yet better than random, policy. An alternative might  
177 be to consider percentile statistics or, if the goal is to evaluate mean performance, to avoid drawing  
178 strong conclusions about individual runs.

179 If the goal is to report mean performance, then a second issue arises. Estimating the mean of these  
180 non-normal performance distributions can be challenging. In Figure 2, approximately 70% of the

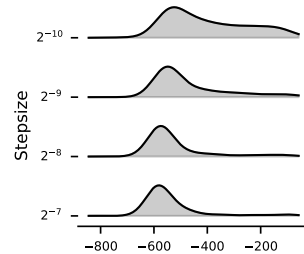
181 density is around a mode centered at 20 return, and the remaining 30% is around a mode centered at  
 182 250 return. As a result, sample means constructed with only three runs are varied and skewed.

183 Further, to report the average performance of the best performing hyperparameter—that is  
 184  $\max_{\theta \in \Theta} \mathbb{E}[G \mid \theta, E]$ —we must first reliably estimate the conditional expected performance for  
 185 each hyperparameter. Computing this expectation can require a large number of samples to obtain a  
 186 reasonable estimate for each hyperparameter. This results in a tradeoff between measuring sensitivity  
 187 and stability: between the breadth of hyperparameter settings that can be studied and the accuracy to  
 188 with which we can feasibly evaluate each hyperparameter.



The summary statistic used to select hyperparameters also interacts with the form of the performance distribution. In the inset figure on the left we show the performance distribution across four different choices of stepsize parameter of DQN in **Cartpole**. If we are interested only in the highest best case performance, then  $2^{-10}$  is preferred. However, if we are particularly concerned with reducing the chances of catastrophic failure (i.e., highest worst case performance), then a stepsize  $2^{-7}$  is preferred. The most common case is to report results for the stepsize with the highest average performance. In this case, a stepsize of  $2^{-9}$  would be preferred.

199 These performance distributions can also look quite different for different  
 200 environments, even with the same algorithm. For Cartpole (above),  
 201 the distribution is increasingly long-tailed with smaller stepsizes. For  
 202 **Puddle World**, shown in the inset figure on the right, the distributions  
 203 are always bimodal with one mode around -600 return and a second  
 204 mode around -200 return. With smaller stepsizes, the density around  
 205 the better performance mode increases, shifting the mean of the distri-  
 206 bution. Peak performance does not change; rather the probability that  
 207 DQN has a good run is higher with small stepsizes. This analysis of  
 208 performance distributions raises an important question: do current RL  
 209 algorithms have consistent hyperparameter settings which perform well across many environments?



## 210 4 The Cross-environment Hyperparameter Setting Benchmark

211 In this section, we describe our new benchmark for evaluating RL algorithm across environments, the  
 212 Cross-environment Hyperparameter Setting Benchmark (CHS). Although it seems natural to evaluate  
 213 across environments, standard empirical practice in RL is not done this way. Understanding across-  
 214 environment sensitivity aligns nicely with the intent of sensitivity analysis: elucidating how well an  
 215 algorithm might perform on new environments without extensive hyperparameter tuning. We argue  
 216 that the CHS 1) better aligns empirical practice with the goals of applied RL, 2) is computationally  
 217 feasible even in complex environments, 3) provides novel insights on old ideas (even with small  
 218 environments), and 4) reduces the chances of accidentally publishing incorrect conclusions due to  
 219 statistical noise.

220 We now reiterate the procedure for the CHS with more details than the high-level procedure given  
 221 in Section 2. The first step (**preliminary sweep**) is to draw a small number of samples  $n_{\text{runs}}$  from  
 222  $\mathbb{P}(G \mid \theta, E)$  for every hyperparameter setting and environment and get the summary estimate  $g(E, \theta)$   
 223 from those samples. Typically, we compute  $g(E, \theta)$  as a sample average to estimate  $\mathbb{E}[N_E(G) \mid E, \theta]$ ,  
 224 where  $N_E : \mathbb{R} \rightarrow \mathbb{R}$  is a **normalization** function that we describe below. Then we aggregate across  
 225 environments to estimate  $g(\theta) \approx \mathbb{E}[\mathbb{E}[N_E(G) \mid E, \theta]]$ , where the outer expectation is with respect  
 226 to environments. Then we **select** a single hyperparameter setting with  $\theta_{\text{CHS}} = \arg \max_{\theta \in \Theta} g(\theta)$ .  
 227 Finally, we draw a large number of samples from  $\mathbb{P}(G \mid \theta_{\text{CHS}}, E)$  for every environment and report  
 228 the same summary statistics  $g(E, \theta_{\text{CHS}})$  and  $g(\theta_{\text{CHS}})$  (**re-evaluation**).

229 In order to compute the expectation over environments we must normalize the performance measures.  
 230 Generally, we cannot expect each environment to produce normalized performance numbers. A  
 231 comprehensive discussion of normalization methods is given in Jordan et al. (2020). We use a lightly  
 232 modified version of the CDF normalization method from Jordan et al. (2020),  $N_E(G) = \text{CDF}(G, E)$ ,  
 233 which is itself an instance of probabilistic performance profiles (Barreto et al., 2010).

234 We first collect the performance of each algorithm, environment, and hyperparameter tuple. Then,  
 235 let  $\mathcal{P}_E$  be the pool of performance statistics  $g$  for every agent—namely a run of each algorithm and  
 236 hyperparameter pair—for a given environment  $E$ . Our goal is to take a given agent’s performance  
 237  $x \in \mathcal{P}_E$  and return a normalized performance. The CDF normalization, for this  $x$  in this environment,  
 238 is

$$\text{CDF}(x, E) = \frac{1}{|\mathcal{P}_E|} \sum_{g \in \mathcal{P}_E} \mathbf{1}(g < x)$$

239 where  $\mathbf{1}$  is the indicator function. This mapping says: what percentage of performance values, across  
 240 all runs for all algorithms and all hyperparameter settings, is lower than my performance  $x$  on this  
 241 particular environment  $E$ ? For example, if  $\text{CDF}(x, E) = 0.25$ , then this agent’s performance is  
 242 quite low in this environment, as only 25% of other agents’ performance was worse and 75% was  
 243 higher, across agents tested. This normalization accounts for the difficulty of the problem, and reflects  
 244 relative performance amongst agents tested. Note that this normalization uses an empirical CDF,  
 245 rather than the true CDF for the environment and set of hyperparameters and agents. This means  
 246 there is a small amount of bias when estimating  $\mathbb{E}[\mathbb{E}[N_E(G) \mid E, \theta]]$ . This bias dissipates with an  
 247 increasing numbers of samples and equally impacts all compared algorithms.

248 Selecting hyperparameters with the CHS can require significantly fewer samples compared with  
 249 conventional per-environment tuning. Per-environment tuning requires a sufficiently accurate estimate  
 250 of the conditional expectation  $\mathbb{E}[G \mid E, \theta]$  for every  $\theta \in \Theta$  and for every  $E \in \mathcal{E}$ , requiring a number  
 251 of runs proportional to  $|\Theta||\mathcal{E}|$ . The CHS, on the other hand, requires only an accurate estimate of  
 252  $\mathbb{E}[N_E(G) \mid \theta] = \mathbb{E}[\mathbb{E}[N_E(G) \mid E, \theta]]$  which requires a number of runs proportional only to  $|\mathcal{E}|$ . By  
 253 designing a process which selects hyperparameters first using a smaller number of runs, we can  
 254 reserve more computational resources for re-evaluation. Once we select the best hyperparameters,  
 255 the cost of collecting samples is independent of  $\Theta$ , and so we can decouple the precision of our  
 256 performance estimate from the number of hyperparameter settings that we evaluate for each algorithm.

257 Finally, we can contrast this benchmark with a recent evaluation scheme that uses random hy-  
 258 perparameter selection (Jordan et al., 2020). In order to capture variation in performance due to  
 259 hyperparameter sensitivity, Jordan et al. (2020) treats hyperparameters as random variables and  
 260 samples according to an experimenter-designated distribution over hyperparameters, reporting the  
 261 mean and uncertainty with respect to this added variance, similar to the procedure used in Jaderberg  
 262 et al. (2016). This evaluation methodology provides some insight into the difficulty of tuning, though  
 263 requires a sensible distribution over hyperparameters to be chosen. The CHS, on the other hand, asks:  
 264 is there a hyperparameter setting for which this algorithm can perform well across environments? It  
 265 motivates instead identifying that single hyperparameter, and potentially fixing it in the algorithm, or  
 266 suggesting that the algorithm needs to be improved so that such a hyperparameter could feasibly be  
 267 found. Both of these strategies help identify algorithms that are difficult to tune, but the CHS is easier  
 268 to use and computationally cheaper.

## 269 5 Evaluating the Cross-environment Hyperparameter Setting Benchmark

270 In this section, we evaluate the CHS by comparing four algorithms across several classic control  
 271 environments. For this evaluation, we require environments where hundreds of independent samples  
 272 of performance can be drawn across a large hyperparameter sweep in a computationally tractable  
 273 way. We emphasize that this is not a general requirement of the CHS and is required only in this  
 274 case of evaluating the CHS’s responsiveness to perturbations in the experimental process. Because  
 275 these classic control environments are cheap to run and provide meaningful insights in differentiating  
 276 modern RL algorithms (Obando-Ceron and Castro, 2021), we name this specific benchmark the  
 277 Small Control CHS (SC-CHS). In Section 6 we provide a realistic demonstration of the CHS on a  
 278 larger dataset with a more complex algorithm.<sup>1</sup>

279 For the following investigations, we compare two deep RL algorithms based on DQN (Mnih et al.,  
 280 2013) and two control algorithms based on linear function approximation using tile-coded features  
 281 (Sutton and Barto, 2018). The deep RL algorithms, DQN and DeepQ, differ only in their loss: DQN  
 282 uses a clipped loss and DeepQ uses a mean squared error. For the two tile-coding agents, QLearning  
 283 is off-policy and bootstraps using the greedy action, while ESARSA is on-policy and bootstraps using  
 284 an expectation over actions. Further details on the algorithms can be found in Appendix C.

<sup>1</sup>All code can be found at <https://github.com/andnp/single-hyperparameter-benchmark>.

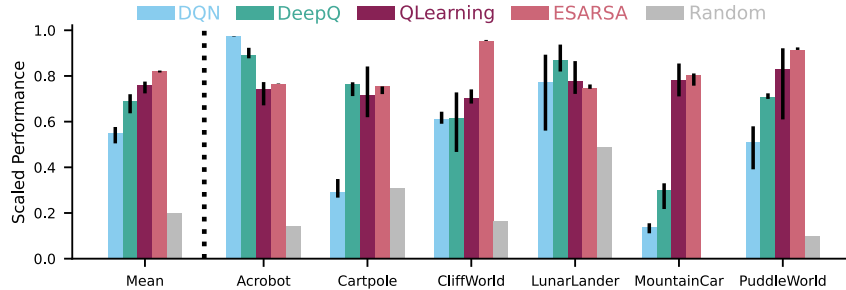


Figure 3: Applying the CHS to 10k simulated experiments. Error bars show 95% bootstrap confidence intervals. Although only three runs were used to select hyperparameters, conclusions about algorithm ranking using the CHS are perfectly consistent across all 10k experiments.

285 The SC-CHS consists of a suite of classic control environments commonly used in RL: Acrobot  
 286 (Sutton, 1996), Cartpole (Barto et al., 1983; Brockman et al., 2016), Cliff World (Sutton and Barto,  
 287 2018), Lunar Lander (Brockman et al., 2016), Mountain Car (Moore, 1990; Sutton, 1996), and Puddle  
 288 World (Sutton, 1996). We used a discount factor of  $\gamma = 0.99$  and a maximum episode length of 500  
 289 steps (except in Cliff World which had a maximum length of 50 steps). We ran all algorithms for  
 290 200k learning steps on each environment except Lunar Lander, where we used 250k learning steps to  
 291 ensure all algorithms have reliably converged. Further details motivating this choice of environments  
 292 can be found in Appendix C.1.

293 We swept over several hyperparameter settings. For all algorithms we swept eight stepsize values,  
 294  $\alpha \in \{2^{-12}, 2^{-11}, \dots, 2^{-5}\}$  for the deep RL algorithms and  $\alpha \in \{2^{-9}, 2^{-8}, \dots, 2^{-2}\}$  for the tile-  
 295 coded algorithms. The deep RL algorithms used experience replay and target networks, so we swept  
 296 over replay buffer sizes of  $\{2000, 4000\}$  and target network refresh rates of  $\{1, 8, 32\}$  steps where a  
 297 one step refresh indicates target networks are not used. The algorithms with tile-coding learn online  
 298 from the most recent sample; we select number of tiles in each tiling in  $\{2, 4, 8\}$  and number of tilings  
 299 in  $\{8, 16, 32\}$ . More details on the other hyperparameters and design decisions are in Appendix C.

300 **Variance over simulated experiments.** Here we demonstrate that the CHS provides low variance  
 301 conclusions over 10k simulated experiments using the benchmark. We use bootstrap sampling  
 302 to compute 10k sample means over three random seeds for every algorithm, environment, and  
 303 hyperparameter to first select hyperparameters using the CHS. We then evaluate the performance  
 304 of each algorithm on each environment with 250 independent runs for the selected hyperparameter  
 305 settings and compare the conclusions for each of the 10k simulated experiments.

306 Figure 3 demonstrates the consistency of conclusions made using the CHS across 10k simulated  
 307 experiments. Using the CHS we would rank algorithms from best to worst ESARSA, QLearning,  
 308 DeepQ, and DQN on this benchmark, and this ranking was successfully detected in every experiment.  
 309 Conclusions on individual environments are less consistent. This is because selecting one hyperpa-  
 310 rameter across all these environments was difficult. In some runs, performance in one environment  
 311 was sacrificed for the performance in the others; in another run, it was a different environment.

312 We provide more insight into the difficulty of selecting a single hyperparameter across problems, in  
 313 Appendix B.1. We additionally show that the distribution of selected hyperparameters with the CHS is  
 314 narrow and consistent over simulated experiments, unlike parameters chosen independently for each  
 315 environment. Because conclusions are often drawn by aggregating results over environments—either  
 316 formally as in the CHS or informally by counting the number of environments where an algorithm  
 317 outperforms others—reporting results over a consistent and narrow distribution of hyperparameters  
 318 leads towards lower variance claims and greater reproducibility. We include results selecting hyperpa-  
 319 rameters according to the worst-case performance across environments in Appendix B.4; the results  
 320 are highly similar, albeit slightly lower variance.

321 The cost of running a single experiment represented in Figure 3 is quite low. The deep RL algorithms  
 322 test 48 hyperparameter settings at a cost of 20 minutes per run, while the tile-coded algorithms test  
 323 72 settings at the cost of two minutes per run. Timings are with respect to a modern 2.1Ghz Intel  
 324 Xeon processor. This comes out to a total of 1762 hours of CPU time to complete three runs for  
 325 hyperparameter selection and 250 runs for evaluation, cheaper than the experiment using 10 runs and  
 326 conventional per-environment tuning shown in Table 1 which cost approximately 2208 hours. The

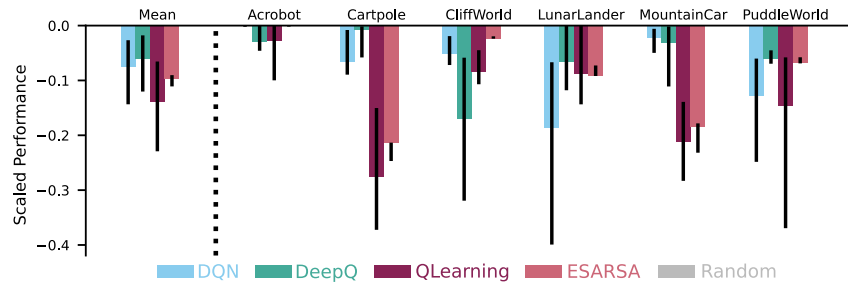


Figure 4: The change in performance for each algorithm on every environment when using the CHS versus conventional per-environment tuning. A larger drop in performance indicates a larger degree of environment overfitting when results are reported with per-environment tuning. Error bars show 95% confidence intervals over 10k bootstrap samples.

327 CHS successfully detected the correct ordering of algorithms in every trial, while the conventional  
 328 per-environment tuning experiment failed to detect the correct ordering with surprising frequency.

329 **The CHS is a less optimistic measure of performance.** A motivating factor for the CHS is providing  
 330 a more challenging benchmark to test across-environment insensitivity to selection of hyperparameters.  
 331 Because algorithms are limited to selecting a single champion hyperparameter setting—as opposed  
 332 to selecting a new hyperparameter setting for every environment—we expect a considerable drop in  
 333 performance under the CHS. We evaluate the extent of this performance drop for our four algorithms  
 334 by first computing near optimal parameters  $\theta^* \in \Theta$  for each environment using the full 250 random  
 335 seeds to obtain high confidence estimates of average performance  $\mathbb{E}[N_E(G) | E, \theta^*]$ . We then apply  
 336 the CHS to select hyperparameters for each algorithm using three random seeds for 10k simulated  
 337 experiments. We report sample estimates of  $\mathbb{E}[N_E(G) | E, \theta^*] - \mathbb{E}[N_E(G) | E, \theta_{\text{CHS}}]$ .

338 In Figure 4 we can see there is substantial drop in reported performance when using the CHS versus  
 339 per-environment tuning. The variance is high, indicating that for some runs, the performance drop  
 340 was substantial: almost 0.4 under our normalization between  $[0,1]$ . Algorithms with a large drop in  
 341 performance indicate more environment-specific overfitting under per-environment tuning. Because  
 342 we swept over many more hyperparameter settings for the tile-coding algorithms than for the deep  
 343 RL algorithms—72 settings versus 48 settings—it is unsurprising that per-environment tuning led to  
 344 far more environment overfitting in the tile-coding algorithms.

345 **Tuning on a subset of environments.** An empirical practice that is highly related to the CHS is using  
 346 a subset of environments to select hyperparameters, then reporting the performance of the selected  
 347 hyperparameters across an entire suite of environments. We refer to this practice as *subset-CHS*.  
 348 This practice is used in the Atari suite for example, where it was suggested to use five of the 57  
 349 games for hyperparameter tuning (Bellemare et al., 2013). To investigate the variance of conclusions  
 350 using the subset-CHS, we run 10k simulated experiments using two of our six environments to select  
 351 hyperparameters. For each of the simulated experiments, we randomly select two environments to  
 352 use for hyperparameter selection. To reduce the variance, we allow each algorithm 100 runs of every  
 353 hyperparameter setting on every environment to perform hyperparameter selection, then evaluate the  
 354 performance on the full 250 runs for the hyperparameter selected by the subset-CHS. More results,  
 355 including with varying number of runs and environments used for hyperparameter selection, can be  
 356 found in Appendix B.

357 In Figure 5, we see that the ordering of algorithms is extremely high-variance—especially compared  
 358 to Figure 3 which uses all six environments to select hyperparameters and only three runs. This  
 359 result also illustrates large differences between individual environments, where the variance on Lunar  
 360 Lander—especially for DQN—suggests that hyperparameters selected for other environments are  
 361 likely to cause worse-than-random performance on Lunar Lander. At least among the four tested  
 362 algorithms, it is clear that hyperparameter sensitivity is too high to use environment subselection to  
 363 reduce the computational burden of hyperparameter tuning.

364 **Bias of the CHS.** Both the CHS and conventional per-environment tuning use biased sample esti-  
 365 mates due to the maximization over hyperparameters. The bias due to maximization over random  
 366 samples is exaggerated both as the set  $\Theta$  grows and as the number of samples used to evaluate  
 367  $\mathbb{E}[G | E, \theta]$  shrinks. We first estimate the true per-environment maximizing parameters  $\theta^*$  and the  
 368 true CHS parameter  $\theta_{\text{CHS}}^*$  using 250 samples for every hyperparameter setting and environment.



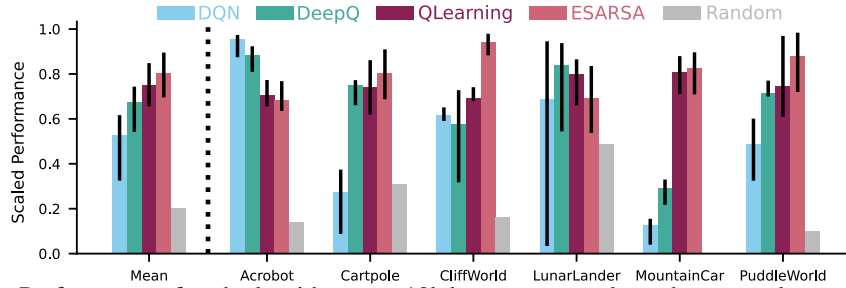


Figure 5: Performance of each algorithm over 10k bootstrap samples, where sample means are computed with 100 runs. Each bootstrap sample randomly selects two environments for hyperparameter tuning, then evaluates the chosen hyperparameter setting on all six environments with 250 runs. Error bars show 95% bootstrap confidence intervals.

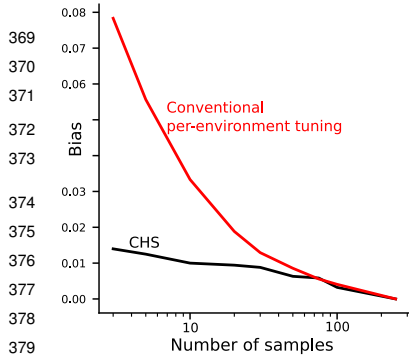


Figure 6: Bias of the CHS vs. per-environment tuning.

We then resample three samples per hyperparameter and environment to simulate an experiment using three seeds to compute sample averages, we select the maximizing parameter of these sample averages,  $\hat{\theta}$ , and we report  $\mathbb{E}[G | E, \theta^*] - \mathbb{E}[G | E, \hat{\theta}]$ . The corresponding procedure is used for the CHS.

In Figure 6, we report the bias of each procedure applied to DQN and the small control domain suite. On the vertical axis we report the bias and on the horizontal axis we show the number of random seeds used to select hyperparameters. As both procedures approach a sufficiently large number of samples to select hyperparameters, the bias of these procedures approaches zero. However when using few random seeds—for instance ten or fewer as is common in the literature—the bias of the conventional method is several times larger than that of the CHS. As a result of this overestimation bias, it is common for results in the literature to present highly optimistic results especially for algorithms with more hyperparameters.

## 6 A Demonstrative Example of Using the CHS

We finish with a large-scale demonstration of our benchmark across the 28 environments of the DMControl suite (Tassa et al., 2018), which we will call the DMC-CHS. For this comparison, we test an open hypothesis in the continuous control literature: does Ornstein-Uhlenbeck (OU) noise (Uhlenbeck and Ornstein, 1930) improve exploration over naive uncorrelated Gaussian noise? Autocorrelated noise for exploration was shown to be beneficial for robotics (Wawrzyński, 2015), inspiring the use of an OU noise process for DDPG (Lillicrap et al., 2016), where a single set of hyperparameters was used across 20 Mujoco environments using five seeds. Later work replaced OU noise with Gaussian noise, noting no difference in performance (Fujimoto et al., 2018; Barth-Maron et al., 2018), but without empirical support for the claim. To the best of our knowledge, no careful empirical investigation of this hypothesis has yet been published.

To apply the DMC-CHS, we first evaluate 36 hyperparameter settings with three runs per environment, for a total of 84 runs to estimate  $\mathbb{E}[N_E(G) | \theta]$  for each  $\theta \in \Theta$ . Then we use 30 runs to evaluate the chosen  $\theta_{\text{CHS}}$  for a total of 840 runs to estimate  $\mathbb{E}[N_E(G) | \theta_{\text{CHS}}]$ . We report the swept hyperparameters as well as the selected  $\theta_{\text{CHS}}$  in Appendix B.5. We use 1k bootstrap samples to compute confidence intervals and report the overall findings in the table in Figure 7. We find that OU noise does not outperform Gaussian noise on the DMC-CHS. Considering even the extremes of the confidence intervals there is no meaningful difference in performance between these exploration methods, suggesting further runs would be unlikely to change our conclusion. We visualize the performance of OU noise on the complete suite, considering Gaussian noise experiments as a baseline in Figure 7. This visualization summarizes whether, and to what degree, OU noise improves upon Gaussian noise in each environment of the DMControl suite. In only 10 of the 28 environments, OU noise improves upon Gaussian noise, with a large improvement only in the *WalkerRun* environment. Additional results are included in Appendix B.5.

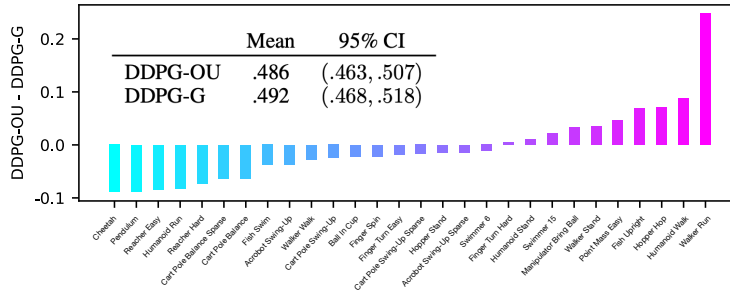


Figure 7: Comparing DDPG using OU noise vs. Gaussian noise across the DMControl suite. The inset table shows the mean performance with 95% confidence interval for the two versions of DDPG used in these experiments. Visualized in the bar plot is the performance of DDPG with OU noise, per environment in the suite, considering DDPG with Gaussian noise as a baseline.

409 **7 Conclusion**

410 In this work, we introduced a new benchmark for evaluating RL algorithms across environments,  
 411 but perhaps more important are the insights we gained. Of the five algorithms we tested (including  
 412 DQN and DDPG), none exhibited good performance on our CHS benchmark; aligning with the  
 413 common view that we do not yet have generally applicable RL algorithms. The CHS benchmark  
 414 produces reliable conclusions with only three runs in the preliminary sweep while providing a new  
 415 challenging aspect to small computationally-cheap environments, allowing small university labs and  
 416 tech giants alike to conduct rigorous and meaningful comparisons. Finally, prior work has disagreed  
 417 on the benefit of using OU or Gaussian noise in DDPG on Mujoco-based environments. Perhaps  
 418 some combination of too few runs, using default hyperparameters, or problematic environment sub-  
 419 selection yielded conflicting results. Our results with CHS suggest there is no significant performance  
 420 difference across a suite of 28 Mujoco environments, putting this debate to bed. The CHS benchmark  
 421 can play a role uncovering falsehoods and resolving disputes.

422 The CHS is a general procedure for evaluating performance across environments. We provide two  
 423 example instantiations of the CHS, the SC-CHS for discrete action control on small domains and  
 424 the DMC-CHS for continuous control on large simulated environments, however the CHS can also  
 425 be extended to use arbitrary environment sets to allow targeted evaluation across environments with  
 426 certain desirable properties. For example, the taxonomies of Atari games identified in Bellemare  
 427 et al. (2016), the off-policy evaluation environments used in Sutton et al. (2009), or the taxonomy  
 428 of exploration environments from Yasui et al. (2019) are each sets of environments that have been  
 429 previously identified and used across the literature. Applying the CHS to any one of the environment  
 430 sets provides a new challenge, and in some small way can push us towards generally applicable RL  
 431 agents.

## 432 **References**

- 433 T. W. Archibald, K. I. M. McKinnon, and L. C. Thomas. On the generation of markov decision  
434 processes. *Journal of the Operational Research Society*, 1995.
- 435 André M.S. Barreto, Heder S. Bernardino, and Helio J.C. Barbosa. Probabilistic performance profiles  
436 for the experimental evaluation of stochastic algorithms. *Conference on Genetic and Evolutionary*  
437 *Computation*, 2010.
- 438 Gabriel Barth-Marón, Matthew W. Hoffman, David Budden, Will Dabney, Dan Horgan, Dhruva TB,  
439 Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributed Distributional Deterministic  
440 Policy Gradients. *International Conference on Learning Representations*, 2018.
- 441 Andrew G. Barto, Richard S. Sutton, and Charles W. Anderson. Neuronlike adaptive elements that can  
442 solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics*,  
443 1983.
- 444 Charles Beattie, Joel Z. Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler,  
445 Andrew Lefrancq, Simon Green, Víctor Valdés, and Amir Sadik. Deepmind lab. *arXiv preprint*  
446 *arXiv:1612.03801*, 2016.
- 447 M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The Arcade Learning Environment: An  
448 Evaluation Platform for General Agents. *Journal of Artificial Intelligence Research*, 2013.
- 449 Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.  
450 Unifying count-based exploration and intrinsic motivation. *Advances in neural information*  
451 *processing systems*, 2016.
- 452 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and  
453 Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 454 Cédric Colas, Olivier Sigaud, and Pierre-Yves Oudeyer. How Many Random Seeds? Statistical  
455 Power Analysis in Deep Reinforcement Learning Experiments. *arXiv:1806.08295*, 2018.
- 456 Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph,  
457 and Aleksander Ma. Implementation Matters In Deep Policy Gradients: A Case Study On PPO  
458 and TRPO. *International Conference on Learning Representations*, 2019.
- 459 Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward, Yotam  
460 Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. IMPALA:  
461 Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures. *Interna-*  
462 *tional Conference on Machine Learning*, 2018.
- 463 Scott Fujimoto, Herke Van Hoof, and David Meger. Addressing function approximation error in  
464 actor-critic methods. *International Conference on Machine Learning*, 2018.
- 465 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.  
466 Deep Reinforcement Learning that Matters. *AAAI*, 2018.
- 467 Matt Hoffman, Bobak Shahriari, John Aslanides, Gabriel Barth-Marón, Feryal Behbahani, Tamara  
468 Norman, Abbas Abdolmaleki, Albin Cassirer, Fan Yang, Kate Baumli, Sarah Henderson, Alex  
469 Novikov, Sergio Gómez Colmenarejo, Serkan Cabi, Caglar Gulcehre, Tom Le Paine, Andrew  
470 Cowie, Ziyu Wang, Bilal Piot, and Nando de Freitas. Acme: A Research Framework for Distributed  
471 Reinforcement Learning. *arXiv:2006.00979*, 2020.
- 472 Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z. Leibo, David  
473 Silver, and Koray Kavukcuoglu. Reinforcement Learning with Unsupervised Auxiliary Tasks.  
474 *International Conference on Learning Representations*, 2016.
- 475 Scott M. Jordan, Yash Chandak, Daniel Cohen, Mengxue Zhang, and Philip S. Thomas. Evaluating  
476 the Performance of Reinforcement Learning Algorithms. *International Conference on Machine*  
477 *Learning*, 2020.
- 478 Diederik P Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *International*  
479 *Conference on Learning Representations*, 2015.

480 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,  
481 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *International Conference on Learning Representations*, 2016.  
482

483 Ashique Rupam Mahmood, Hado Van Hasselt, and Richard S. Sutton. Weighted importance sampling  
484 for off-policy learning with linear function approximation. *Advances in Neural Information*  
485 *Processing Systems*, 2014.

486 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
487 Wierstra, and Martin Riedmiller. Playing Atari with Deep Reinforcement Learning. 2013.

488 Andrew William Moore. *Efficient Memory-Based Learning for Robot Control*. PhD thesis, University  
489 of Cambridge, 1990.

490 Johan S. Obando-Ceron and Pablo Samuel Castro. Revisiting Rainbow: Promoting more insightful  
491 and inclusive deep reinforcement learning research. *International Conference on Machine Learning*,  
492 2021.

493 Andrew Patterson, Adam White, Sina Ghiassian, and Martha White. A Generalized Projected  
494 Bellman Error for Off-policy Value Estimation in Reinforcement Learning. *In submission*, 2021.

495 Joelle Pineau, Philippe Vincent-Lamarre, Koustuv Sinha, Vincent Larivière, Alina Beygelzimer,  
496 Florence d’Alché-Buc, Emily Fox, and Hugo Larochelle. Improving Reproducibility in Machine  
497 Learning Research (A Report from the NeurIPS 2019 Reproducibility Program). *arXiv:2003.12206*,  
498 2020.

499 Antonin Raffin, Ashley Hill, Maximilian Ernestus, Adam Gleave, Anssi Kanervisto, and Noah  
500 Dormann. Stable Baselines3. DLR-RM, 2019.

501 Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized Experience Replay.  
502 *International Conference on Learning Representations*, 2016.

503 Harm Seijen and Rich Sutton. True online TD ( $\lambda$ ). In *International Conference on Machine*  
504 *Learning*, 2014.

505 Richard S. Sutton. Generalization in reinforcement learning: Successful examples using sparse coarse  
506 coding. *Advances in Neural Information Processing Systems*, 1996.

507 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.

508 Richard S. Sutton, Hamid Reza Maei, Doina Precup, Shalabh Bhatnagar, David Silver, Csaba  
509 Szepesvári, and Eric Wiewiora. Fast gradient-descent methods for temporal-difference learning  
510 with linear function approximation. *International Conference on Machine Learning*, 2009.

511 Adrien Ali Taïga, William Fedus, Marlos C. Machado, Aaron Courville, and Marc G. Bellemare.  
512 Benchmarking Bonus-Based Exploration Methods on the Arcade Learning Environment. *Exploration in Reinforcement Learning Workshop, ICML*, 2019.  
513

514 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden,  
515 Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Riedmiller.  
516 DeepMind Control Suite. *arXiv:1801.00690*, 2018.

517 George E. Uhlenbeck and Leonard S. Ornstein. On the theory of the Brownian motion. *Physical*  
518 *review*, 1930.

519 Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-  
520 learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.

521 Hado van Hasselt, Matteo Hessel, and John Aslanides. When to use parametric models in reinforce-  
522 ment learning? *arXiv:1906.05243*, 2019.

523 Ziyu Wang, Tom Schaul, Matteo Hessel, Hado van Hasselt, Marc Lanctot, and Nando de Freitas.  
524 Dueling Network Architectures for Deep Reinforcement Learning. *International Conference on*  
525 *Machine Learning*, 2016.

- 526 Paweł Wawrzyński. Control Policy with Autocorrelated Noise in Reinforcement Learning for  
527 Robotics. *International Journal of Machine Learning and Computing*, 2015.
- 528 Adam White and Martha White. Investigating practical linear temporal difference learning. *International Conference on Autonomous Agents and Multi-Agent Systems*, 2016.
- 530 Shimon Whiteson, Brian Tanner, Matthew E. Taylor, and Peter Stone. Generalized domains for  
531 empirical evaluations in reinforcement learning. *Workshop on Evaluation Methods for Machine  
532 Learning at ICML*, 2009.
- 533 Niko Yasui, Sungsu Lim, Cam Linke, Adam White, and Martha White. An Empirical and Conceptual  
534 Categorization of Value-based Exploration Methods. *Exploration in Reinforcement Learning  
535 Workshop, ICML*, 2019.

536 **Checklist**

- 537 1. For all authors...
- 538 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
539 contributions and scope? [Yes]
- 540 (b) Did you describe the limitations of your work? [Yes] Section 4 discusses our relation-  
541 ship to Jordan et al. (2020) and that their work allows further quantitative evaluation of  
542 the difficulty in choosing hyperparameters.
- 543 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 544 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
545 them? [Yes]
- 546 2. If you are including theoretical results...
- 547 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 548 (b) Did you include complete proofs of all theoretical results? [N/A]
- 549 3. If you ran experiments (e.g. for benchmarks)...
- 550 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
551 mental results (either in the supplemental material or as a URL)? [Yes]
- 552 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
553 were chosen)? [Yes]
- 554 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
555 ments multiple times)? [Yes]
- 556 (d) Did you include the total amount of compute and the type of resources used (e.g.,  
557 type of GPUs, internal cluster, or cloud provider)? [Yes] Appendix A details the  
558 computational resources required for this paper.
- 559 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 560 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 561 (b) Did you mention the license of the assets? [N/A]
- 562 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 563
- 564 (d) Did you discuss whether and how consent was obtained from people whose data you're  
565 using/curating? [N/A]
- 566 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
567 information or offensive content? [N/A]
- 568 5. If you used crowdsourcing or conducted research with human subjects...
- 569 (a) Did you include the full text of instructions given to participants and screenshots, if  
570 applicable? [N/A]
- 571 (b) Did you describe any potential participant risks, with links to Institutional Review  
572 Board (IRB) approvals, if applicable? [N/A]
- 573 (c) Did you include the estimated hourly wage paid to participants and the total amount  
574 spent on participant compensation? [N/A]