
OPEN RL BENCHMARK: Comprehensive Tracked Experiments for Reinforcement Learning

Shengyi Huang^{1,2*} Quentin Gallouédec^{1,3*} Florian Felten⁴ Antonin Raffin⁵
Rousslan Fernand Julien Dossa⁶ Yanxiao Zhao^{7,8} Ryan Sullivan⁹ Viktor Makoviychuk¹⁰
Denys Makoviichuk¹¹ Mohamad H. Danesh¹² Cyril Rourégous¹³ Jiayi Weng
Chufan Chen¹⁴ Md Masudur Rahman¹⁵ João G. M. Araújo¹⁶ Guorui Quan¹⁷
Daniel C.H. Tan^{18,19} Timo Klein^{20,21} Rujikorn Charakorn²² Mark Towers²³
Yann Berthelot^{24,25} Kinal Mehta²⁶ Dipam Chakraborty²⁷ Arjun KG
Valentin Charrat²⁸ Chang Ye²⁹ Zichen Liu³⁰ Lucas N. Alegre³¹ Alexander Nikulin³²
Xiao Hu³³ Tianlin Liu³⁴ Jongwook Choi³⁵ Brent Yi³⁶

Abstract

1 In many Reinforcement Learning (RL) papers, learning curves are useful indicators
2 to measure the effectiveness of RL algorithms. However, the complete raw data
3 of the learning curves are rarely available. As a result, it is usually necessary
4 to reproduce the experiments from scratch, which can be time-consuming and
5 error-prone. We present OPEN RL BENCHMARK (ORLB), a set of fully tracked
6 RL experiments, including not only the usual data such as episodic return, but also
7 all algorithm-specific and system metrics. ORLB is community-driven: anyone
8 can download, use, and contribute to the data. At the time of writing, more than
9 25,000 runs have been tracked, for a cumulative duration of more than 8 years.
10 It covers a wide range of RL libraries and reference implementations. Special
11 care is taken to ensure that each experiment is precisely reproducible by providing
12 not only the full parameters, but also the versions of the dependencies used to
13 generate it. In addition, ORLB comes with a command-line interface (CLI) for
14 easy fetching and generating figures to present the results. In this document, we
15 include two case studies to demonstrate the usefulness of ORLB in practice. To
16 the best of our knowledge, ORLB is the first RL benchmark of its kind, and the
17 authors hope that it will improve and facilitate the work of researchers in the field.

18 1 Introduction

19 Reinforcement Learning (RL) research is based on comparing new methods to baselines to assess
20 progress (Patterson et al., 2023). This process requires the availability of the data associated with
21 these baselines (Raffin et al., 2021) or, alternatively, the ability to replicate them and generate the
22 data oneself (Raffin, 2020). In addition, reproducible results allow the methods to be compared with
23 new benchmarks and to identify the areas in which the methods excel and those in which they are
24 likely to fail, thus providing avenues for future research.

25 In practice, the RL research community faces complex challenges in comparing new methods with
26 reference data. The unavailability of reference data requires researchers to reproduce experiments,
27 which is difficult due to insufficient source code documentation and evolving software dependencies.

*Equal contributions

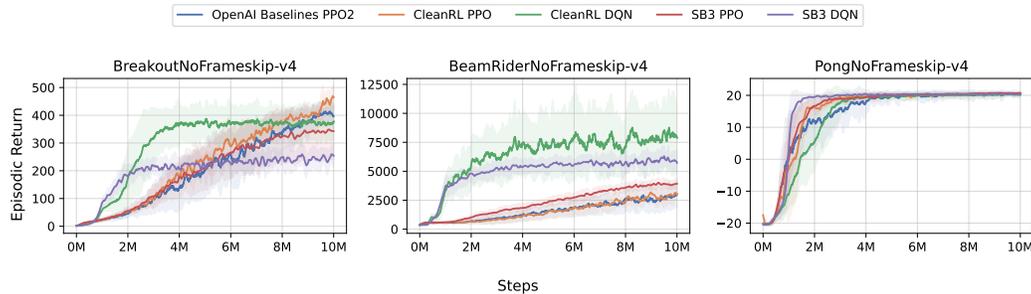


Figure 1: Example of learning curves obtained with OPEN RL BENCHMARK. These compare the episodic returns obtained by different implementations of PPO and DQN on three Atari games.

28 Implementation details, as highlighted in past research, can significantly impact results (Henderson
 29 et al., 2018; Huang et al., 2022a). Moreover, limited computing resources play a crucial role, hindering
 30 the reproduction process and affecting researchers without substantial access.

31 The lack of standardized metrics and benchmarks across studies not only impedes comparison but
 32 also results in a substantial waste of time and resources. To address these issues, the RL community
 33 must establish rigorous reproducibility standards, ensuring replicability and comparability across
 34 studies. Transparent sharing of data, code, and experimental details, along with the adoption of
 35 consistent metrics and benchmarks, would collectively enhance the evaluation and progression of RL
 36 research, ultimately accelerating advancements in the field.

37 ORLB presents a rich collection of tracked RL experiments and aims to set a new standard by
 38 providing a diverse training dataset. This initiative prioritizes the use of existing data over re-running
 39 baselines, emphasizing reproducibility and transparency. Our contributions are:

- 40 • **Extensive dataset:** Offers a large, diverse collection of tracked RL experiments.
- 41 • **Standardization:** Establishes a new norm by encouraging reliance on existing data, reducing
 42 the need for re-running baselines.
- 43 • **Comprehensive metrics:** Includes diverse tracked metrics for method-specific and system
 44 evaluation, in addition to episodic return.
- 45 • **Reproducibility:** Emphasizes clear instructions and fixed dependencies, ensuring easy
 46 experiment replication.
- 47 • **Resource for research:** Serves as a valuable and collaborative resource for RL research.
- 48 • **Facilitating exploration:** Enables reliable exploration and assessment of new and existing
 49 RL methods.

50 2 Comprehensive overview of ORLB: content, methodology, tools, and 51 applications

52 This section provides a detailed exploration of the contents of ORLB, including its diverse set of
 53 libraries and environments, and the metrics it contains. We also look at the practical aspects of using
 54 ORLB, highlighting its ability to ensure accurate reproducibility and facilitate the creation of data
 55 visualizations thanks to its CLI.

56 2.1 Content

57 ORLB data is stored and shared with Weights and Biases (Biewald, 2020). The data is contained
 58 in a common entity named `openrlbenchmark`. Runs are divided into several *projects*. A project
 59 can correspond to a library, but it can also correspond to a set of more specific runs, such as
 60 `envpool-cleanrl` in which we find CleanRL runs (Huang et al., 2022b) launched with the EnvPool

61 implementation of environments (Weng et al., 2022b). A project can also correspond to a reference
 62 implementation, such as TD3 (project sfujim-TD3) or Phasic Policy Gradient (Cobbe et al., 2021)
 63 (project phasic-policy-gradient). ORLB also includes reports, which are interactive documents
 64 designed to enhance the visualization of selected representations. These reports provide a more
 65 user-friendly format for practitioners to share, discuss, and analyze experimental results, even across
 66 different projects. Figure 2 shows a preview of one such report.

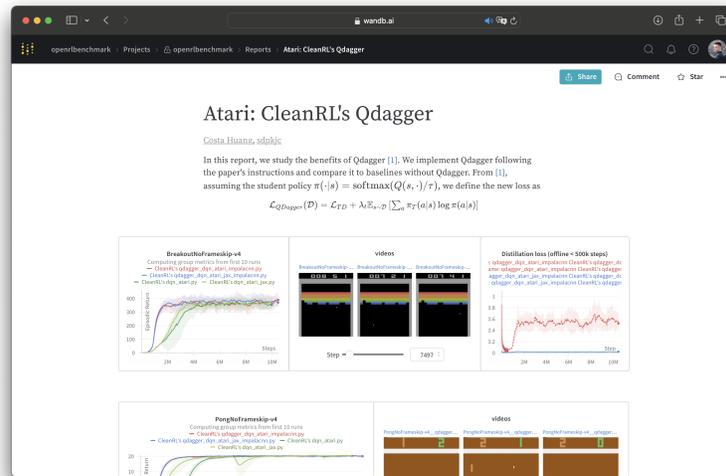


Figure 2: An example of a report on the Weights and Biases platform, dealing with the contribution of QDagger (Agarwal et al., 2022), and using data from ORLB. The URL to access the report is <https://wandb.ai/openrlbenchmark/openrlbenchmark/reports/Atari-CleanRL-s-Qdagger--Vml1dz0ONTg1ODY5>.

67 At the time of writing, ORLB contains nearly 25,000 runs, for a total of 72,000 hours (more than 8
 68 years) of tracking. In the following paragraphs, we present the libraries and environments for which
 69 runs are available in ORLB, as well as the metrics tracked.

70 **Libraries** ORLB contains runs for several reference RL libraries. These libraries are: abcdRL
 71 (Zhao, 2022), Acme (Hoffman et al., 2020), Cleanba (Huang et al., 2023), CleanRL (Huang et al.,
 72 2022b), jaxrl (Kostrikov, 2021), moolib (Mella et al., 2022), MORL-Baselines (Felten et al., 2023),
 73 OpenAI Baselines (Dhariwal et al., 2017), rlgames (Makoviichuk & Makoviychuk, 2021) Stable
 74 Baselines3 (Raffin et al., 2021; Raffin, 2020) Stable Baselines Jax (Raffin et al., 2021) and TorchBeast
 75 (Küttler et al., 2019).

76 **Environments** The runs contained in ORLB cover a wide range of classic environments. They
 77 include Atari (Bellemare et al., 2013; Machado et al., 2018), Classic control (Brockman et al., 2016),
 78 Box2d (Brockman et al., 2016) and MuJoCo (Todorov et al., 2012) as part of either Gym (Brockman
 79 et al., 2016) or Gymnasium (Towers et al., 2023) or EnvPool (Weng et al., 2022b). They also
 80 include Bullet (Coumans & Bai, 2016), Procgen Benchmark (Cobbe et al., 2020), Fetch environments
 81 (Plappert et al., 2018), PandaGym (Gallouédec et al., 2021), highway-env (Leurent, 2018), Minigrid
 82 (Chevalier-Boisvert et al., 2023) and MO-Gymnasium (Alegre et al., 2022).

83 **Tracked metrics** Metrics are recorded throughout the learning process, consistently linked with a
 84 global step indicating the number of interactions with the environment, and an absolute time, which
 85 allows to compute the duration of a run. We categorize these metrics into four distinct groups:

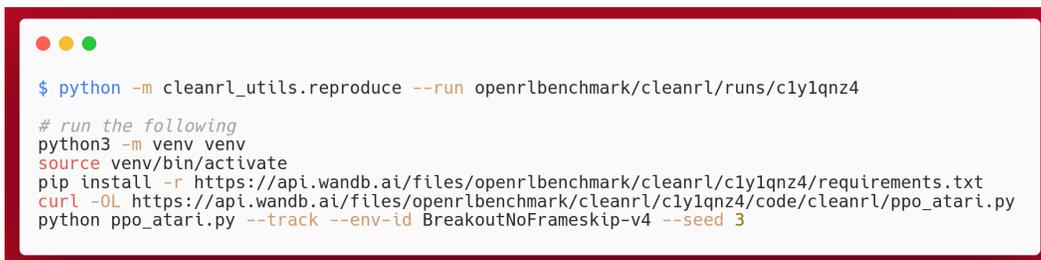
- 86 • **Training-related metrics:** These are general metrics related to RL learning. This category
 87 contains, for example, the average returns obtained, the episode length or the number of
 88 collected samples per second.

- 89 • **Method-specific metrics:** These are losses and measures of key internal values of the
90 methods. For PPO, for example, this category includes the value loss, the policy loss, the
91 entropy or the approximate KL divergence.
- 92 • **Evolving configuration parameters:** These are configuration values that change during the
93 learning process. This category includes, for example, the learning rate when there is decay,
94 or the exploration rate (ϵ) in the Deep Q-Network (DQN) (Mnih et al., 2013).
- 95 • **System metrics:** These are metrics related to system components. These could be GPU
96 memory usage, its power consumption, its temperature, system and process memory usage,
97 CPU usage or even network traffic.

98 The specific metrics available may vary from one library to another. In addition, even where the
99 metrics are technically similar, the terminology or key used to record them may vary from one
100 library to another. Users are advised to consult the documentation specific to each library for precise
101 information on these measures.

102 2.2 Everything you need for perfect repeatability

103 Reproducing experimental results in computational research, as discussed in Section 4.3, is often
104 challenging due to evolving codebases, incomplete hyperparameter listings, version discrepancies,
105 and compatibility issues. Our approach aims to enhance reproducibility by ensuring users can
106 exactly replicate benchmark results. Each experiment includes a complete configuration with all
107 hyperparameters, frozen versions of dependencies, and the exact command, including the necessary
108 random seed, for systematic reproducibility. As an example, CleanRL (Huang et al., 2022b) introduces
109 a unique utility that streamlines the process of experiment replication (see Figure 3). This tool
110 produces the command lines to set up a Python environment with the necessary dependencies,
111 download the run file, and the precise command required for the experiment reproduction. Such
112 an approach to reproduction facilitates research and makes it possible to study in depth unusual
113 phenomena, or cases of rupture², in learning processes, which are generally ignored in the results
114 presented, either because they are deliberately left out or because they are erased by the averaging
115 process.



```

$ python -m cleanrl_utils.reproduce --run openrlbenchmark/cleanrl/runs/c1y1qnz4

# run the following
python3 -m venv venv
source venv/bin/activate
pip install -r https://api.wandb.ai/files/openrlbenchmark/cleanrl/c1y1qnz4/requirements.txt
curl -OL https://api.wandb.ai/files/openrlbenchmark/cleanrl/c1y1qnz4/code/cleanrl/ppo_atari.py
python ppo_atari.py --track --env-id BreakoutNoFrameskip-v4 --seed 3

```

Figure 3: CleanRL’s module `reproduce` allows the user to generate, from an ORLB run reference, the exact command suite for an identical reproduction of the run.

116 2.3 The CLI for generating figures in one command line

117 ORLB offers convenient access to raw data from RL libraries on standard environments. It includes a
118 feature for easily extracting and visualizing data in a paper-friendly format, streamlining the process
119 of filtering and extracting relevant runs and metrics for research papers through a single command.
120 The CLI is a powerful tool for generating most metrics-related figures for RL research and notably,
121 all figures in this document were generated using the CLI. The data in ORLB can also be accessed
122 by custom scripts, as detailed in Appendix A.2. Specifically, the CLI integrated into ORLB provides
123 users with the flexibility to:

²Exemplified in <https://github.com/DLR-RM/r1-baselines3-zoo/issues/427>

- 124 • Specify algorithms’ implementations (from which library) along with their corresponding
125 git commit or tag;
- 126 • Choose target environments for analysis;
- 127 • Define the metrics of interest;
- 128 • Opt for the additional generation of metrics and plots using RLiab (Agarwal et al., 2021).

129 Concrete example usage of the CLI and resulting plots are available in Appendix A.1.

130 3 ORLB in action: an insight into case studies

131 ORLB offers a powerful tool for researchers to evaluate and compare different RL algorithms. In this
132 section, we will explore two case studies that showcase its benefits. First, we propose to investigate
133 the effect of using TD(λ) for value estimation in PPO (Schulman et al., 2017) versus using Monte
134 Carlo (MC). This simple study illustrates the use of ORLB through a classic research question.
135 Moreover, to the best of our knowledge, this question has never been studied in the literature. We
136 then show how ORLB is used to demonstrate the speedup and variance reduction of a new IMPALA
137 implementation proposed by Huang et al. (2023). By using ORLB, we can save time and resources
138 while ensuring consistent and reproducible comparisons. These case studies highlight the role of the
139 benchmark in providing insights that can advance the field of RL research.

140 3.1 Easily assess the contribution of TD(λ) for value estimation in PPO

141 In the first case study, we show how ORLB can be used to easily compare the performance of
142 different methods for estimating the value function in PPO (Schulman et al., 2017), one of the
143 many implementation details of this algorithm (Huang et al., 2022a). Specifically, we compare the
144 commonly used Temporal Difference (TD)(λ) estimate to the Monte-Carlo (MC) estimate.

145 PPO typically employs Generalized Advantage Estimation (GAE) (Schulman et al., 2016) to update
146 the actor. The advantage estimate is expressed as follows:

$$A_t^{\text{GAE}(\gamma, \lambda)} = \sum_{l=0}^{N-1} (\gamma \lambda)^l \delta_{t+l}^V \quad (1)$$

147 where $\lambda \in [0, 1]$ adjusts the bias-variance tradeoff and $\delta_{t+l}^V = R_{t+l} + \gamma \hat{V}(S_{t+l+1}) - \hat{V}(S_{t+l})$. The
148 target return for critic optimization is estimated with TD(λ) as follows:

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+n} \quad (2)$$

149 where $G_{t:t+n} = \sum_{k=0}^{n-1} \gamma^k R_{t+k+1} + \gamma^n V(S_{t+n})$ is the n -steps return. In practice, the target return
150 for updating the critic is computed from the GAE value, by adding the minibatch return, a detail
151 usually overlooked by practitioners (Huang et al., 2022a, point 5). While previous studies (Patterson
152 et al., 2023) have shown the joint benefit of GAE and over MC estimates for actor and critic, we
153 focus on the value function alone. To isolate the influence of the value function estimation, we vary
154 the method used for the value function and keep GAE for advantage estimation.

155 The first step is to identify the reference runs in ORLB. Since PPO is a well-known baseline, there
156 are many runs available; we decided to use those from Stable Baselines3 for this example. We
157 then retrieve the exact source code and command used to generate the runs – thanks to the pinned
158 dependencies that come with them – and make the necessary changes to the source code. For each
159 selected environment, we start three learning runs using the same command as the one we retrieved.
160 The runs are saved in a dedicated project³. For fast and user-friendly rendering of the results, we

³<https://wandb.ai/modanesh/openrlbenchmark>

161 create a Weights and Biases report⁴. Using ORLB CLI, we generate Figure 4 and 5. The command
 162 used to generate the figures is given in Appendix B.

163 Figures 4 and 5 give an overview of the results, while detailed plots in the Appendix B provide a
 164 closer look at each environment. The proposed modification to the PPO value function estimation
 165 has an impact on the performance for Atari games (Figure 4a): not using TD(λ) results in lower
 166 scores. However, PPO with MC estimates has similar performance to the original PPO in Box2D
 167 and MuJoCo environments. This example shows how ORLB can be used to quickly investigate the
 168 influence of design choices in RL. It provides baseline results and tools to compare and reproduce
 169 results.

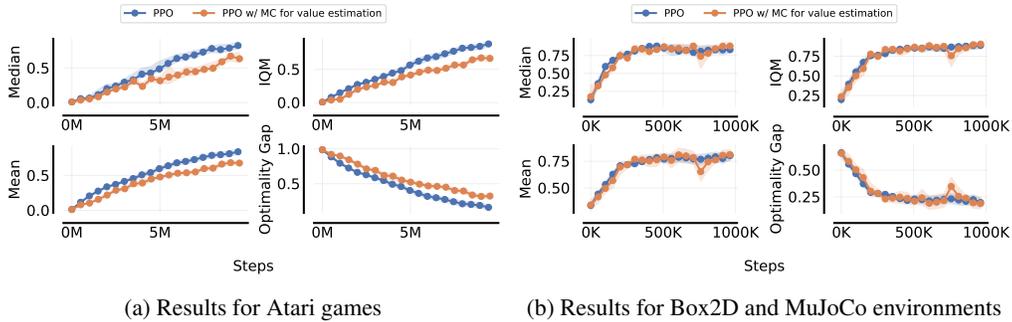


Figure 4: Comparing the original PPO and the PPO with Monte-Carlo (MC) for value estimation. These experiments were conducted over 15 environments, including Atari games, Box2D, and MuJoCo. The plot shows min-max normalized scores with 95% stratified bootstrap CIs.

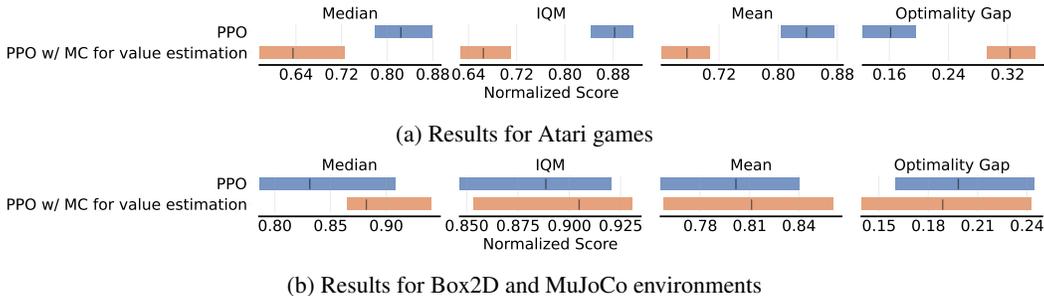


Figure 5: Study of the contribution of GAE for estimating the value used to update the critic in PPO, compared against its variant which uses the MC estimator instead. Figures show the aggregated min-max normalized scores with stratified 95% stratified bootstrap CIs.

170 3.2 Demonstrating the utility of ORLB through the Cleanba case study

171 This section describes how ORLB was instrumental in the evaluation and presentation of Cleanba
 172 (Huang et al., 2023), a new open-source platform for distributed RL implementing highly optimized
 173 distributed variants of PPO (Schulman et al., 2017) and IMPALA (Espeholt et al., 2018). Cleanba’s
 174 authors asserted three points: (1) Cleanba implementations compare favorably with baselines in terms
 175 of sample efficiency, (2) for the same system, the Cleanba implementation is more optimized and
 176 therefore faster, and (3) the design choices allow a reduction in the variability of results.

177 To prove these assertions, the evaluation of Cleanba encountered a common problem in RL research:
 178 the works that initially proposed these baselines did not provide the raw results of their experiments.
 179 Although a reference implementation is available⁵, it is no longer maintained. Subsequent works
 180 such as Moolib (Mella et al., 2022) and TorchBeast (Küttler et al., 2019) have successfully replicated

⁴<https://api.wandb.ai/links/modanesh/izf4yje4>

⁵https://github.com/google-deepmind/scalable_agent

181 the IMPALA results. However, these shared results are limited to the paper’s presented curves, which
 182 provide a smoothed measure of episodic return as a function of interaction steps on a specific set of
 183 Atari tasks. It is worth noting that these tasks are not an exact match for the widely recognized Atari
 184 57, and the raw data used to generate these curves is unavailable.

185 Recognizing the lack of raw data for existing IMPALA implementations, the authors reproduced the
 186 experiments, tracked the runs and integrated them into ORLB. As a reminder, these logged data
 187 include not only the return curves, but also the system configurations and temporal data, which are
 188 crucial to support the Cleanba authors’ optimization claim. Comparable experiments have been run,
 189 tracked and shared on ORLB with the proposed Cleanba implementation.

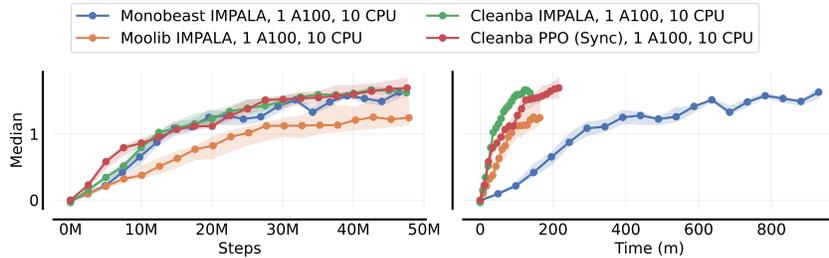


Figure 6: Median human-normalized scores with 95% stratified bootstrap CIs of Cleanba (Huang et al., 2023) variants compared with moolib (Mella et al., 2022) and monobeast (Küttler et al., 2019). The experiments were conducted on 57 Atari games (Bellemare et al., 2013). The data used to generate the figure comes from ORLB, and the figure was generated with a single command from ORLB’s CLI. Figure from (Huang et al., 2023).

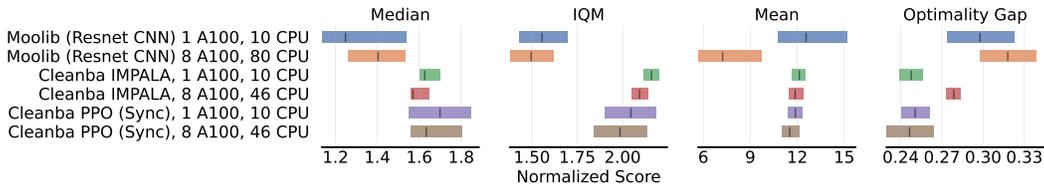


Figure 7: Aggregated normalized human scores with stratified 95% bootstrap CIs, showing that unlike moolib (Mella et al., 2022), Cleanba (Huang et al., 2023) variants have more predictable learning curves (using the same hyperparameters) across different hardware configurations. Figure from (Huang et al., 2023).

190 Using ORLB CLI, the authors generated several figures. In Figure 6, taken from (Huang et al.,
 191 2023), the authors show that the results in terms of sample efficiency compare favorably with the
 192 baselines, and that for the same system configuration, convergence was temporally faster with the
 193 proposed implementation, thus proving claims (1) and (2). Figure 7 demonstrates that Cleanba
 194 variants maintain consistent learning curves across different hardware configurations. Conversely,
 195 moolib’s IMPALA shows marked variability in similar settings, despite identical hyperparameters,
 196 confirming the authors’ third claim.

197 4 Current practices in RL: data reporting, sharing and reproducibility

198 Many new methods have emerged in recent years, with some becoming standard baselines, but
 199 current practices in the field make it challenging to interpret, compare, and replicate study results. In
 200 this section, we highlight the inconsistent presentation of results, focusing on learning curves as an
 201 example. This inconsistency can hinder interpretation and lead to incorrect conclusions. We also note
 202 the insufficient availability of learning data, despite some positive efforts, and examine challenges
 203 related to method reproducibility.

204 4.1 Analyzing learning curve practices

205 Plotting learning curves is a common way to show an agent’s performance over learning. We closely
206 examine the components of learning curves and the choices made by key publications. We find a lack
207 of uniformity, with presentation choices rarely explained and sometimes not explicitly stated.

208 **Axis** Typically, the y axis measures either the return acquired during data collection or evaluation.
209 Some older papers, like (Schulman et al., 2015; Mnih et al., 2016; Schulman et al., 2017), fail to
210 specify the metric, using the vague term *learning curve*. The first approach sums the rewards collected
211 during agent rollout (Dabney et al., 2018; Burda et al., 2019). The second approach suspends training,
212 averaging the agent’s return over episodes, deactivating exploration elements (Fujimoto et al., 2018;
213 Haarnoja et al., 2018; Hessel et al., 2018; Janner et al., 2019; Badia et al., 2020b; Ecoffet et al., 2021;
214 Chen et al., 2021). This method is prevalent and provides a more precise evaluation. Regarding the x
215 axis, while older baselines (Schulman et al., 2015; Mnih et al., 2016) use policy updates and learning
216 epochs, the norm is to use interaction counts with the environment. In Atari environments, it is often
217 the number of frames, adjusting for frame skipping to match human interaction frequency.

218 **Shaded area** Data variability is typically shown with a shaded area, but its definition varies across
219 studies. Commonly, it represents the standard deviation (Chen et al., 2021; Janner et al., 2019) and
220 less commonly half the standard deviation (Fujimoto et al., 2018). Haarnoja et al. (2018) uses a
221 min-max representation to include outliers, covering the entire observed range. This method offers
222 a comprehensive view but amplifies outliers’ impact with more runs. Ecoffet et al. (2021) adopts
223 a probabilistic approach, showing a 95% bootstrap confidence interval around the mean, ensuring
224 statistical confidence. Unfortunately, Schulman et al. (2015, 2017); Mnih et al. (2016); Dabney et al.
225 (2018); Badia et al. (2020b) omit statistical details or even the shaded area, introducing uncertainty in
226 data variability interpretation, as seen in (Hessel et al., 2018).

227 **Normalization and aggregation** Performance aggregation assesses method results across various
228 tasks and domains, indicating their generality and robustness. Outside the Atari context, aggregation
229 practices are uncommon due to the lack of a universal normalization standard. Without a widely
230 accepted normalization strategy, scores are typically not aggregated, or if they are, it relies on a
231 min-max approach lacking absolute significance and unsuitable for comparisons. Early Atari research
232 did not use normalization or aggregate results (Mnih et al., 2013). There has been a shift towards
233 normalizing against human performance, though this has weaknesses and may not reflect true agent
234 mastery (Toromanoff et al., 2019). Aggregation methods vary: the mean is common but influenced
235 by outliers, leading some studies to prefer the more robust median, as in (Hessel et al., 2018). Many
236 papers now report both mean and median results (Dabney et al., 2018; Hafner et al., 2023; Badia
237 et al., 2020a). Recent approaches, like using the Interquartile Mean (IQM), provide a more accurate
238 performance representation across diverse games (Lee et al., 2022), as suggested by Agarwal et al.
239 (2021).

240 4.2 Spectrum of data sharing practices

241 While the mentioned studies often have reference implementations (see Section 4.3), the sharing of
242 training data typically extends only to the curves presented in their articles. This necessitates reliance
243 on libraries that replicate these methods, offering benchmarks with varying levels of completeness.
244 Several widely-used libraries in the field provide high-level summaries or graphical representations
245 without including raw data (e.g., Tensorforce (Kuhnle et al., 2017), Garage (garage contributors,
246 2019), ACME (Hoffman et al., 2020), MushroomRL (D’Eramo et al., 2021), ChainerRL (Fujita
247 et al., 2021), and TorchRL (Bou et al., 2023)). Spinning Up (Achiam, 2018) offers partial data
248 accessibility, providing benchmark curves but withholding raw data. TF-Agent (Guadarrama et al.,
249 2018) is slightly better, offering experiment tracking with links to TensorBoard.dev, though its future
250 is uncertain due to service closure. Tianshou (Weng et al., 2022a) provides individual run reward data
251 for Atari and average rewards for MuJoCo, with more detailed MuJoCo data available via a Google

252 Drive link, but it is not widely promoted. RLLib (Liang et al., 2018) maintains an intermediate
253 stance in data sharing, hosting run data in a dedicated repository. However, this data is specific to
254 select experiments and often presented in non-standard, undocumented formats, complicating its
255 use. Leading effective data-sharing platforms include Dopamine (Castro et al., 2018) and Sample
256 Factory (Petrenko et al., 2020). Dopamine consistently provides accessible raw evaluation data for
257 various seeds and visualizations, along with trained agents on Google Cloud. Sample Factory offers
258 comprehensive data via Weights and Biases (Biewald, 2020) and a selection of pre-trained agents on
259 the Hugging Face Hub, enhancing reproducibility and collaborative research efforts.

260 4.3 Review on reproducibility

261 The literature shows variations in these practices. Some older publications like (Schulman et al.,
262 2015, 2017; Bellemare et al., 2013; Mnih et al., 2016; Hessel et al., 2018) and even recent ones
263 like (Reed et al., 2022) lack a codebase but provide detailed descriptions for replication⁶. However,
264 challenges arise because certain hyperparameters, important but often unreported, can significantly
265 affect performance (Andrychowicz et al., 2020). In addition, implementation choices have proven to
266 be critical (Henderson et al., 2018; Huang et al., 2023, 2022a; Engstrom et al., 2020), complicating
267 the distinction between implementation-based improvements and methodological advances.

268 Recognizing these challenges, the RL community is advocating for higher standards. NeurIPS, for
269 instance, has been requesting a reproduction checklist since 2019 (Pineau et al., 2021). Recent
270 efforts focus on systematic sharing of source code to promote reproducibility. However, codebases
271 are often left unmaintained post-publication (with rare exceptions (Fujimoto et al., 2018)), creating
272 complexity for users dealing with various dependencies and unsolved issues. To address these
273 challenges, libraries have aggregated multiple baseline implementations (see Section 2.1), aiming
274 to match reported paper performance. However, long-term sustainability remains a concern. While
275 these libraries enhance reproducibility, in-depth repeatability is still rare.

276 5 Discussion and conclusion

277 Reproducing results in RL research is often difficult due to limited access to data and code, as well
278 as the impact of minor implementation variations on performance. Researchers typically rely on
279 imprecise comparisons with paper figures, making the reproduction process time-consuming and
280 challenging. To address these issues, we introduce ORLB, a large collection of tracked experiments
281 spanning various algorithms, libraries and benchmarks. ORLB records all relevant metrics and
282 data points, offering detailed resources for precise reproduction. This tool facilitates access to
283 comprehensive datasets, simplifies the extraction of valuable information, enables metric comparisons,
284 and provides a CLI for easier data access and visualization. As a dynamic resource, ORLB is regularly
285 updated by both its maintainers and the user community, gradually improving the reliability of the
286 available results.

287 Despite its strengths, ORLB faces challenges in user-friendliness that need to be addressed. Incon-
288 sistencies between libraries in evaluation strategies and terminology can make it difficult for users.
289 Scaling community engagement becomes a challenge with more members, libraries, and runs. The
290 lack of Git-like version tracking for runs adds to these limitations.

291 ORLB is a major step forward in addressing the needs of RL research. It offers a comprehensive,
292 accessible, and collaborative experiment database, enabling precise comparisons and analysis. It
293 improves data access and promotes a deeper understanding of algorithmic performance. While
294 challenges remain, ORLB has the potential to raise the standard of RL research.

⁶This section uses the taxonomy introduced by Lynnerup et al. (2019): *repeatability* means accurately duplicating an experiment with source code and random seed availability, *reproducibility* involves redoing an experiment using an existing codebase, and *replicability* aims to achieve similar results independently through algorithm implementation.

295 **Affiliations**

- 296 ¹Hugging Face
297 ²Drexel University
298 ³Univ. Lyon, Centrale Lyon, CNRS, INSA Lyon, UCBL, LIRIS, UMR 5205
299 ⁴SnT, University of Luxembourg
300 ⁵German Aerospace Center (DLR) RMC, Weßling, Germany
301 ⁶Graduate School of System Informatics, Kobe University, Hyogo, Japan
302 ⁷School of Computer Science and Technology, University of Chinese Academy of Sciences
303 ⁸Chengdu Institute of Computer Applications, Chinese Academy of Sciences
304 ⁹University of Maryland, College Park
305 ¹⁰NVIDIA
306 ¹¹Snap Inc.
307 ¹²School of Computer Science, McGill University
308 ¹³Polytech Montpellier DO
309 ¹⁴Zhejiang University
310 ¹⁵Department of Computer Science, Purdue University
311 ¹⁶Work done while at Cohere
312 ¹⁷Chinese University of Hong Kong, Shenzhen
313 ¹⁸University College London
314 ¹⁹Agency for Science, Technology and Research
315 ²⁰Faculty of Computer Science, University of Vienna, Vienna, Austria
316 ²¹UniVie Doctoral School Computer Science, University of Vienna
317 ²²Vidyasirimedhi Institute of Science and Technology (VISTEC)
318 ²³University of Southampton
319 ²⁴Univ. Lille, Inria, CNRS, Centrale Lille, UMR 9189 – CRIStAL
320 ²⁵Saint-Gobain Research Paris
321 ²⁶International Institute of Information Technology, Hyderabad, India
322 ²⁷Aicrowd SA
323 ²⁸Valeo Driving Assistance Research
324 ²⁹New York University
325 ³⁰Sea AI Lab
326 ³¹Institute of Informatics, Federal University of Rio Grande do Sul
327 ³²AIRI
328 ³³Department of Automation, Tsinghua University
329 ³⁴University of Basel
330 ³⁵University of Michigan
331 ³⁶UC Berkley

332 **Acknowledgments**

333 This work has been supported by a highly committed RL community. We have listed all the
334 contributors to date, and would like to thank all future contributors and users in advance.

335 This work was granted access to the HPC resources of IDRIS under the allocation 2022-
336 [AD011012172R1] made by GENCI. The MORL-Baselines experiments have been conducted on the
337 HPCs of the University of Luxembourg, and of the Vrije Universiteit Brussel. This work was partly
338 supported by the National Key Research and Development Program of China (2023YFB3308601),
339 Science and Technology Service Network Initiative (KFJ-ST-S-QYZD-2021-21-001), the Talents by
340 Sichuan provincial Party Committee Organization Department, and Chengdu - Chinese Academy
341 of Sciences Science and Technology Cooperation Fund Project (Major Scientific and Technological
342 Innovation Projects). Some experiments are conducted at Stability AI and Hugging Face’s cluster.

343 References

- 344 Joshua Achiam. Spinning Up in Deep Reinforcement Learning. [https://github.com/openai/](https://github.com/openai/spinningup)
345 spinningup, 2018. URL <https://github.com/openai/spinningup>.
- 346 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C. Courville, and Marc G. Bellemare.
347 Deep Reinforcement Learning at the Edge of the Statistical Precipice. In Marc’ Aurelio Ranzato,
348 Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Ad-*
349 *vances in Neural Information Processing Systems 34: Annual Conference on Neural Information*
350 *Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 29304–29320, 2021.
- 351 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C. Courville, and Marc G. Belle-
352 mare. Reincarnating Reinforcement Learning: Reusing Prior Computation to Accelerate Progress.
353 In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.),
354 *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Infor-*
355 *mation Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - De-*
356 *cember 9, 2022*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/ba1c5356d9164bb64c446a4b690226b0-Abstract-Conference.html)
357 [ba1c5356d9164bb64c446a4b690226b0-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/ba1c5356d9164bb64c446a4b690226b0-Abstract-Conference.html).
- 358 Lucas N. Alegre, Florian Felten, El-Ghazali Talbi, Grégoire Danoy, Ann Nowé, Ana L. C. Bazzan, and
359 Bruno C. da Silva. MO-Gym: A Library of Multi-Objective Reinforcement Learning Environments.
360 In *Proceedings of the 34th Benelux Conference on Artificial Intelligence BNAIC/Benelearn 2022*,
361 2022.
- 362 Marcin Andrychowicz, Anton Raichuk, Piotr Stanczyk, Manu Orsini, Sertan Girgin, Raphaël Marinier,
363 Léonard Hussenot, Matthieu Geist, Olivier Pietquin, Marcin Michalski, Sylvain Gelly, and Olivier
364 Bachem. What Matters In On-Policy Reinforcement Learning? A Large-Scale Empirical Study.
365 *arXiv preprint arXiv:2006.05990*, 2020.
- 366 Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi,
367 Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the Atari Human Benchmark.
368 In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18*
369 *July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 507–517.
370 PMLR, 2020a. URL <http://proceedings.mlr.press/v119/badia20a.html>.
- 371 Adrià Puigdomènech Badia, Pablo Sprechmann, Alex Vitvitskyi, Zhaohan Daniel Guo, Bilal Piot,
372 Steven Kapturowski, Olivier Tieleman, Martín Arjovsky, Alexander Pritzel, Andrew Bolt, and
373 Charles Blundell. Never Give Up: Learning Directed Exploration Strategies. In *8th International*
374 *Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.
375 OpenReview.net, 2020b. URL <https://openreview.net/forum?id=Sye57xStvB>.
- 376 Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The Arcade Learning
377 Environment: An Evaluation Platform for General Agents. *Journal of Artificial Intelligence*
378 *Research*, 47:253–279, 2013. doi: 10.1613/JAIR.3912. URL [https://doi.org/10.1613/](https://doi.org/10.1613/jair.3912)
379 [jair.3912](https://doi.org/10.1613/jair.3912).
- 380 Lukas Biewald. Experiment Tracking with Weights and Biases, 2020. URL [https://www.wandb.](https://www.wandb.com/)
381 [com/](https://www.wandb.com/). Software available from wandb.com.
- 382 Albert Bou, Matteo Bettini, Sebastian Dittert, Vikash Kumar, Shagun Sodhani, Xiaomeng Yang,
383 Gianni De Fabritiis, and Vincent Moens. TorchRL: A Data-Driven Decision-Making Library for
384 Pytorch. *arXiv preprint arXiv:2306.00577*, 2023.
- 385 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and
386 Wojciech Zaremba. OpenAI Gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 387 Yuri Burda, Harrison Edwards, Amos J. Storkey, and Oleg Klimov. Exploration by random network
388 distillation. In *7th International Conference on Learning Representations, ICLR 2019, New*
389 *Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL [https://openreview.net/](https://openreview.net/forum?id=H1lJJnR5Ym)
390 [forum?id=H1lJJnR5Ym](https://openreview.net/forum?id=H1lJJnR5Ym).

- 391 Pablo Samuel Castro, Subhodeep Moitra, Carles Gelada, Saurabh Kumar, and Marc G. Belle-
392 mare. Dopamine: A Research Framework for Deep Reinforcement Learning. *arXiv preprint*
393 *arXiv:1812.06110*, 2018.
- 394 Xinyue Chen, Che Wang, Zijian Zhou, and Keith W. Ross. Randomized Ensembled Double Q-
395 Learning: Learning Fast Without a Model. In *9th International Conference on Learning Rep-*
396 *resentations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL
397 <https://openreview.net/forum?id=AY8zfZmOtDd>.
- 398 Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem
399 Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & Miniworld: Modular &
400 Customizable Reinforcement Learning Environments for Goal-Oriented Tasks. *arXiv preprint*
401 *arXiv:2306.13831*, 2023.
- 402 Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging Procedural Generation
403 to Benchmark Reinforcement Learning. In *Proceedings of the 37th International Conference on*
404 *Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of*
405 *Machine Learning Research*, pp. 2048–2056. PMLR, 2020. URL [http://proceedings.mlr.](http://proceedings.mlr.press/v119/cobbe20a.html)
406 [press/v119/cobbe20a.html](http://proceedings.mlr.press/v119/cobbe20a.html).
- 407 Karl Cobbe, Jacob Hilton, Oleg Klimov, and John Schulman. Phasic Policy Gradient. In Marina
408 Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine*
409 *Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine*
410 *Learning Research*, pp. 2020–2027. PMLR, 2021. URL [http://proceedings.mlr.press/](http://proceedings.mlr.press/v139/cobbe21a.html)
411 [v139/cobbe21a.html](http://proceedings.mlr.press/v139/cobbe21a.html).
- 412 Erwin Coumans and Yunfei Bai. PyBullet, a Python Module for Physics Simulation for Games,
413 Robotics and Machine Learning. 2016.
- 414 Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit Quantile Networks for
415 Distributional Reinforcement Learning. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings*
416 *of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan,*
417 *Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp.
418 1104–1113. PMLR, 2018. URL <http://proceedings.mlr.press/v80/dabney18a.html>.
- 419 Carlo D’Eramo, Davide Tateo, Andrea Bonarini, Marcello Restelli, and Jan Peters. MushroomRL:
420 Simplifying Reinforcement Learning Research. *Journal of Machine Learning Research*, 22(131):
421 1–5, 2021. URL <http://jmlr.org/papers/v22/18-056.html>.
- 422 Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford,
423 John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. OpenAI Baselines. [https://](https://github.com/openai/baselines)
424 github.com/openai/baselines, 2017. URL <https://github.com/openai/baselines>.
- 425 Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O. Stanley, and Jeff Clune. First Return,
426 Then Explore. *Nature*, 590(7847):580–586, 2021. doi: 10.1038/S41586-020-03157-9. URL
427 <https://doi.org/10.1038/s41586-020-03157-9>.
- 428 Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph,
429 and Aleksander Madry. Implementation Matters in Deep RL: A Case Study on PPO and
430 TRPO. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa,*
431 *Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL [https://openreview.net/forum?](https://openreview.net/forum?id=r1etN1rtPB)
432 [id=r1etN1rtPB](https://openreview.net/forum?id=r1etN1rtPB).
- 433 Lasse Espeholt, Hubert Soyer, Rémi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward, Yotam
434 Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. IMPALA:
435 Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures. In Jen-
436 nifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on*

- 437 *Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, vol-
438 ume 80 of *Proceedings of Machine Learning Research*, pp. 1406–1415. PMLR, 2018. URL
439 <http://proceedings.mlr.press/v80/espeholt18a.html>.
- 440 Florian Felten, Lucas Nunes Alegre, Ann Nowe, Ana L. C. Bazzan, El Ghazali Talbi, Grégoire
441 Danoy, and Bruno Castro da Silva. A Toolkit for Reliable Benchmarking and Research in Multi-
442 Objective Reinforcement Learning. In *Proceedings of the Neural Information Processing Systems*
443 *Track on Datasets and Benchmarks 3, NeurIPS Datasets and Benchmarks 2023*, 2023. URL
444 <https://openreview.net/forum?id=jfwRLudQyj>.
- 445 Scott Fujimoto, Herke van Hoof, and David Meger. Addressing Function Approximation Error
446 in Actor-Critic Methods. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the*
447 *35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm,*
448 *Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1582–
449 1591. PMLR, 2018. URL <http://proceedings.mlr.press/v80/fujimoto18a.html>.
- 450 Yasuhiro Fujita, Prabhat Nagarajan, Toshiki Kataoka, and Takahiro Ishikawa. ChainerRL: A Deep
451 Reinforcement Learning Library. *Journal of Machine Learning Research*, 22(77):1–14, 2021.
452 URL <http://jmlr.org/papers/v22/20-376.html>.
- 453 Quentin Gallouédec, Nicolas Cazin, Emmanuel Dellandréa, and Liming Chen. panda-gym: Open-
454 Source Goal-Conditioned Environments for Robotic Learning. *4th Robot Learning Workshop:*
455 *Self-Supervised and Lifelong Learning at NeurIPS*, 2021.
- 456 The garage contributors. Garage: A toolkit for reproducible reinforcement learning research. <https://github.com/rlworkgroup/garage>, 2019.
- 457
- 458 Sergio Guadarrama, Anoop Korattikara, Oscar Ramirez, Pablo Castro, Ethan Holly, Sam Fishman,
459 Ke Wang, Ekaterina Gonina, Neal Wu, Efi Kokipoulou, Luciano Sbaiz, Jamie Smith, Gábor
460 Bartók, Jesse Berent, Chris Harris, Vincent Vanhoucke, and Eugene Brevdo. TF-Agents: A library
461 for Reinforcement Learning in TensorFlow. <https://github.com/tensorflow/agents>, 2018.
462 URL <https://github.com/tensorflow/agents>.
- 463 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft Actor-Critic: Off-Policy
464 Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In Jennifer G. Dy and
465 Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning,*
466 *ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings*
467 *of Machine Learning Research*, pp. 1856–1865. PMLR, 2018. URL <http://proceedings.mlr.press/v80/haarnoja18b.html>.
- 468
- 469 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering Diverse Domains
470 through World Models. *arXiv preprint arXiv:2301.04104*, 2023.
- 471 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.
472 Deep Reinforcement Learning That Matters. In Sheila A. McIlraith and Kilian Q. Weinberger
473 (eds.), *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18),*
474 *the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium*
475 *on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA,*
476 *February 2-7, 2018*, pp. 3207–3214. AAAI Press, 2018. doi: 10.1609/AAAI.V32I1.11694. URL
477 <https://doi.org/10.1609/aaai.v32i1.11694>.
- 478 Matteo Hessel, Joseph Modayil, Hado van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney,
479 Dan Horgan, Bilal Piot, Mohammad Gheshlaghi Azar, and David Silver. Rainbow: Combining
480 Improvements in Deep Reinforcement Learning. In Sheila A. McIlraith and Kilian Q. Weinberger
481 (eds.), *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18),*
482 *the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium*
483 *on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA,*
484 *February 2-7, 2018*, pp. 3215–3222. AAAI Press, 2018. doi: 10.1609/AAAI.V32I1.11796. URL
485 <https://doi.org/10.1609/aaai.v32i1.11796>.

486 Matthew W. Hoffman, Bobak Shahriari, John Aslanides, Gabriel Barth-Maron, Nikola Momchev,
487 Danila Sinopalnikov, Piotr Stańczyk, Sabela Ramos, Anton Raichuk, Damien Vincent, Léonard
488 Hussenot, Robert Dadashi, Gabriel Dulac-Arnold, Manu Orsini, Alexis Jacq, Johan Ferret, Nino
489 Vieillard, Seyed Kamyar Seyed Ghasemipour, Sertan Girgin, Olivier Pietquin, Feryal Behbahani,
490 Tamara Norman, Abbas Abdolmaleki, Albin Cassirer, Fan Yang, Kate Baumli, Sarah Henderson,
491 Abe Friesen, Ruba Haroun, Alex Novikov, Sergio Gómez Colmenarejo, Serkan Cabi, Caglar
492 Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Andrew Cowie, Ziyu Wang, Bilal Piot, and Nando
493 de Freitas. Acme: A Research Framework for Distributed Reinforcement Learning. *arXiv preprint*
494 *arXiv:2006.00979*, 2020.

495 Shengyi Huang, Rousslan Fernand Julien Dossa, Antonin Raffin, Anssi Kanervisto, and
496 Weixun Wang. The 37 Implementation Details of Proximal Policy Optimization.
497 In *ICLR Blog Track*, 2022a. URL [https://iclr-blog-track.github.io/2022/](https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/)
498 [03/25/ppo-implementation-details/](https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/). [https://iclr-blog-track.github.io/2022/03/25/ppo-](https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/)
499 [implementation-details/](https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/).

500 Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Kinal
501 Mehta, and João G.M. Araújo. CleanRL: High-quality Single-file Implementations of Deep
502 Reinforcement Learning Algorithms. *Journal of Machine Learning Research*, 23(274):1–18,
503 2022b. URL <http://jmlr.org/papers/v23/21-1342.html>.

504 Shengyi Huang, Jiayi Weng, Rujikorn Charakorn, Min Lin, Zhongwen Xu, and Santiago Ontañón.
505 Cleanba: A Reproducible and Efficient Distributed Reinforcement Learning Platform, 2023.

506 Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to Trust Your Model:
507 Model-Based Policy Optimization. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelz-
508 imer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neu-*
509 *ral Information Processing Systems 32: Annual Conference on Neural Information Pro-*
510 *cessing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*,
511 pp. 12498–12509, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/5f4f461eff3099671ad63c6f3f094f7f-Abstract.html)
512 [5f4f461eff3099671ad63c6f3f094f7f-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/5f4f461eff3099671ad63c6f3f094f7f-Abstract.html).

513 Ilya Kostrikov. JAXRL: Implementations of Reinforcement Learning algorithms in JAX. <https://github.com/ikostrikov/jaxrl>,
514 Oct 2021. URL [https://github.com/ikostrikov/](https://github.com/ikostrikov/jaxrl)
515 [jaxrl](https://github.com/ikostrikov/jaxrl).

516 Alexander Kuhnle, Michael Schaarschmidt, and Kai Fricke. Tensorforce: a TensorFlow library
517 for applied reinforcement learning. <https://github.com/tensorforce/tensorforce>, 2017.
518 URL <https://github.com/tensorforce/tensorforce>.

519 Heinrich Küttler, Nantas Nardelli, Thibaut Lavril, Marco Selvatici, Viswanath Sivakumar, Tim
520 Rocktäschel, and Edward Grefenstette. TorchBeast: A PyTorch Platform for Distributed RL. *arXiv*
521 *preprint arXiv:1910.03552*, 2019.

522 Kuang-Huei Lee, Ofir Nachum, Mengjiao Yang, Lisa Lee, Daniel Freeman, Sergio Guadarrama, Ian
523 Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, and Igor Mordatch. Multi-Game Decision
524 Transformers. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh
525 (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural*
526 *Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 -*
527 *December 9, 2022*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/b2cac94f82928a85055987d9fd44753f-Abstract-Conference.html)
528 [b2cac94f82928a85055987d9fd44753f-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/b2cac94f82928a85055987d9fd44753f-Abstract-Conference.html).

529 Edouard Leurent. An Environment for Autonomous Driving Decision-Making. <https://github.com/eleurent/highway-env>, 2018. URL <https://github.com/eleurent/highway-env>.

531 Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph
532 Gonzalez, Michael I. Jordan, and Ion Stoica. RLlib: Abstractions for Distributed Reinforcement
533 Learning. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International*

534 *Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15,*
535 *2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 3059–3068. PMLR, 2018.
536 URL <http://proceedings.mlr.press/v80/liang18b.html>.

537 Nicolai A. Lynnerup, Laura Nolling, Rasmus Hasle, and John Hallam. A Survey on Repro-
538 ducibility by Evaluating Deep Reinforcement Learning Algorithms on Real-World Robots. In
539 Leslie Pack Kaelbling, Danica Kragic, and Komei Sugiura (eds.), *3rd Annual Conference on*
540 *Robot Learning, CoRL 2019, Osaka, Japan, October 30 - November 1, 2019, Proceedings*, vol-
541 ume 100 of *Proceedings of Machine Learning Research*, pp. 466–489. PMLR, 2019. URL
542 <http://proceedings.mlr.press/v100/lynnrup20a.html>.

543 Marlos C. Machado, Marc G. Bellemare, Erik Talvitie, Joel Veness, Matthew J. Hausknecht, and
544 Michael Bowling. Revisiting the Arcade Learning Environment: Evaluation Protocols and Open
545 Problems for General Agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018. doi:
546 10.1613/JAIR.5699. URL <https://doi.org/10.1613/jair.5699>.

547 Denys Makoviichuk and Viktor Makoviychuk. rl-games: A High-performance Framework for
548 Reinforcement Learning. https://github.com/Denys88/rl_games, May 2021. URL https://github.com/Denys88/rl_games.

550 Vegard Mella, Eric Hambro, Danielle Rothermel, and Heinrich Küttler. moolib: A Platform for
551 Distributed RL. *GitHub repository*, 2022. URL [https://github.com/facebookresearch/](https://github.com/facebookresearch/moolib)
552 [moolib](https://github.com/facebookresearch/moolib).

553 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
554 Wierstra, and Martin A. Riedmiller. Playing Atari with Deep Reinforcement Learning. *arXiv*
555 *preprint arXiv:1312.5602*, 2013.

556 Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim
557 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous Methods for Deep Reinforcement
558 Learning. In Maria-Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of the 33rd*
559 *International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24,*
560 *2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pp. 1928–1937. JMLR.org,
561 2016. URL <http://proceedings.mlr.press/v48/mniha16.html>.

562 Andrew Patterson, Samuel Neumann, Martha White, and Adam White. Empirical Design in Rein-
563 forcement Learning. *arXiv preprint arXiv:2304.01315*, 2023.

564 Aleksei Petrenko, Zhehui Huang, Tushar Kumar, Gaurav S. Sukhatme, and Vladlen Koltun. Sample
565 Factory: Egocentric 3D Control from Pixels at 100000 FPS with Asynchronous Reinforcement
566 Learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML*
567 *2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*,
568 pp. 7652–7662. PMLR, 2020. URL [http://proceedings.mlr.press/v119/petrenko20a.](http://proceedings.mlr.press/v119/petrenko20a.html)
569 [html](http://proceedings.mlr.press/v119/petrenko20a.html).

570 Joelle Pineau, Philippe Vincent-Lamarre, Koustuv Sinha, Vincent Larivière, Alina Beygelzimer,
571 Florence d’Alché-Buc, Emily B. Fox, and Hugo Larochelle. Improving Reproducibility in Machine
572 Learning Research (A Report from the NeurIPS 2019 Reproducibility Program). *Journal of*
573 *Machine Learning Research*, 22:164:1–164:20, 2021. URL [http://jmlr.org/papers/v22/](http://jmlr.org/papers/v22/20-303.html)
574 [20-303.html](http://jmlr.org/papers/v22/20-303.html).

575 Matthias Plappert, Marcin Andrychowicz, Alex Ray, Bob McGrew, Bowen Baker, Glenn Powell,
576 Jonas Schneider, Josh Tobin, Maciek Chociej, Peter Welinder, Vikash Kumar, and Wojciech
577 Zaremba. Multi-Goal Reinforcement Learning: Challenging Robotics Environments and Request
578 for Research. *arXiv preprint arXiv:1802.09464*, 2018.

579 Antonin Raffin. RL Baselines3 Zoo. <https://github.com/DLR-RM/rl-baselines3-zoo>, 2020.

- 580 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
581 Dormann. Stable-Baselines3: Reliable Reinforcement Learning Implementations. *Journal of*
582 *Machine Learning Research*, 22(268):1–8, 2021.
- 583 Scott E. Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov,
584 Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom
585 Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol
586 Vinyals, Mahyar Bordbar, and Nando de Freitas. A Generalist Agent. *Transactions on Machine*
587 *Learning Research*, 2022, 2022. URL <https://openreview.net/forum?id=1ikK0kHjvj>.
- 588 John Schulman, Sergey Levine, Pieter Abbeel, Michael I. Jordan, and Philipp Moritz. Trust Region
589 Policy Optimization. In Francis R. Bach and David M. Blei (eds.), *Proceedings of the 32nd*
590 *International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*,
591 volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 1889–1897. JMLR.org, 2015.
592 URL <http://proceedings.mlr.press/v37/schulman15.html>.
- 593 John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan, and Pieter Abbeel. High-
594 Dimensional Continuous Control Using Generalized Advantage Estimation. In Yoshua Ben-
595 gio and Yann LeCun (eds.), *4th International Conference on Learning Representations, ICLR*
596 *2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL
597 <http://arxiv.org/abs/1506.02438>.
- 598 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy
599 Optimization Algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 600 Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A Physics Engine for Model-Based
601 Control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS*
602 *2012, Vilamoura, Algarve, Portugal, October 7-12, 2012*, pp. 5026–5033. IEEE, 2012.
- 603 Marin Toromanoff, Émilie Wirbel, and Fabien Moutarde. Is Deep Reinforcement Learning Really
604 Superhuman on Atari? *arXiv preprint arXiv:1908.04683*, 2019.
- 605 Mark Towers, Jordan K. Terry, Ariel Kwiatkowski, John U. Balis, Gianluca de Cola, Tristan Deleu,
606 Manuel Goulão, Andreas Kallinteris, Arjun KG, Markus Krimmel, Rodrigo Perez-Vicente, Andrea
607 Pierré, Sander Schulhoff, Jun Jet Tai, Andrew Tan Jin Shen, and Omar G. Younis. Gymnasium,
608 March 2023. URL <https://zenodo.org/record/8127025>.
- 609 Jiayi Weng, Huayu Chen, Dong Yan, Kaichao You, Alexis Duburcq, Minghao Zhang, Yi Su, Hang Su,
610 and Jun Zhu. Tianshou: A Highly Modularized Deep Reinforcement Learning Library. *Journal*
611 *of Machine Learning Research*, 23(267):1–6, 2022a. URL [http://jmlr.org/papers/v23/](http://jmlr.org/papers/v23/21-1127.html)
612 [21-1127.html](http://jmlr.org/papers/v23/21-1127.html).
- 613 Jiayi Weng, Min Lin, Shengyi Huang, Bo Liu, Denys Makoviichuk, Viktor Makoviychuk, Zichen
614 Liu, Yufan Song, Ting Luo, Yukun Jiang, Zhongwen Xu, and Shuicheng Yan. EnvPool: A Highly
615 Parallel Reinforcement Learning Environment Execution Engine. In *Proceedings of the Neural*
616 *Information Processing Systems Track on Datasets and Benchmarks 2, NeurIPS Datasets and*
617 *Benchmarks 2022*, 2022b. URL [http://papers.nips.cc/paper_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/8caaf08e49ddb6694fae067442ee21-Abstract-Datasets_and_Benchmarks.html)
618 [8caaf08e49ddb6694fae067442ee21-Abstract-Datasets_and_Benchmarks.html](http://papers.nips.cc/paper_files/paper/2022/hash/8caaf08e49ddb6694fae067442ee21-Abstract-Datasets_and_Benchmarks.html).
- 619 Yanxiao Zhao. abcdRL: Modular Single-file Reinforcement Learning Algorithms Library. [https:](https://github.com/sdpkjc/abcdrl)
620 [//github.com/sdpkjc/abcdrl](https://github.com/sdpkjc/abcdrl), December 2022. URL [https://github.com/sdpkjc/](https://github.com/sdpkjc/abcdrl)
621 [abcdrl](https://github.com/sdpkjc/abcdrl).

622 Checklist

- 623 1. For all authors...

- 624 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
625 contributions and scope? [Yes]
- 626 (b) Did you describe the limitations of your work? [Yes] see Section 5.
- 627 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 628 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
629 them? [Yes]
- 630 2. If you are including theoretical results...
- 631 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 632 (b) Did you include complete proofs of all theoretical results? [N/A]
- 633 3. If you ran experiments (e.g. for benchmarks)...
- 634 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
635 mental results (either in the supplemental material or as a URL)? [Yes] The paper deals
636 specifically with new ways of sharing experimental results to improve reproducibility.
- 637 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
638 were chosen)? [Yes]
- 639 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
640 ments multiple times)? [Yes]
- 641 (d) Did you include the total amount of compute and the type of resources used (e.g., type
642 of GPUs, internal cluster, or cloud provider)? [Yes]
- 643 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 644 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 645 (b) Did you mention the license of the assets? [Yes]
- 646 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
647 Each experiment on ORLB is carefully linked to the source code needed to produce it.
- 648 (d) Did you discuss whether and how consent was obtained from people whose data you're
649 using/curating? [N/A] Data is collected by proactive contributors
- 650 (e) Did you discuss whether the data you are using/curating contains personally identifiable
651 information or offensive content? [N/A]
- 652 5. If you used crowdsourcing or conducted research with human subjects...
- 653 (a) Did you include the full text of instructions given to participants and screenshots, if
654 applicable? [N/A]
- 655 (b) Did you describe any potential participant risks, with links to Institutional Review
656 Board (IRB) approvals, if applicable? [N/A]
- 657 (c) Did you include the estimated hourly wage paid to participants and the total amount
658 spent on participant compensation? [N/A]

659 A Plotting results guidelines

660 A.1 Using the CLI

661 This section gives notable additional examples of usage of the provided CLI. A more comprehensive
662 set of examples and manual is available in the README page of the project.

663 A.1.1 Plotting episodic return from various libraries

664 First, we showcase the most basic usage of the CLI, that is comparing two different implementations
665 of the same algorithm based on learning curve of episodic return. For example, Figure 8 and 9
666 compare CleanRL’s TD3 implementation against the original TD3, both in terms of sample efficiency and
667 time. The command used to generate this plot is listed below.

```
668 python -m openrlbenchmark.rlops \  
669     --filters '?we=openrlbenchmark&wpn=sfujim-TD3&ceik=env&cen=policy&metric=charts/episodic_return' 'TD3?  
670     cl=Official TD3' \  
671     --filters '?we=openrlbenchmark&wpn=cleanrl&ceik=env_id&cen=exp_name&metric=charts/episodic_return' '  
672     td3_continuous_action_jax?cl=Clean RL TD3' \  
673     --env-ids HalfCheetah-v2 Walker2d-v2 Hopper-v2 \  
674     --pc.ncols 3 \  
675     --pc.ncols-legend 2 \  
676     --output-filename static/td3_vs_cleanrl \  
677     --scan-history
```

678 In the above command, `wpn` denotes the project name, typically the learning library name. This allows
679 to fetch results of implementations from different projects. Moreover, it is possible to specify which
680 metric to compare, in this case `charts/episodic_return`. Also, the CLI provides the possibility
681 to select a given algorithm and apply a different name in the plot, e.g. we rename TD3 to *Official*
682 *TD3* and `td3_continuous_action_jax` to *Clean RL TD3*. Finally, we can also select a set of
683 environments through the `--end-ids` option.

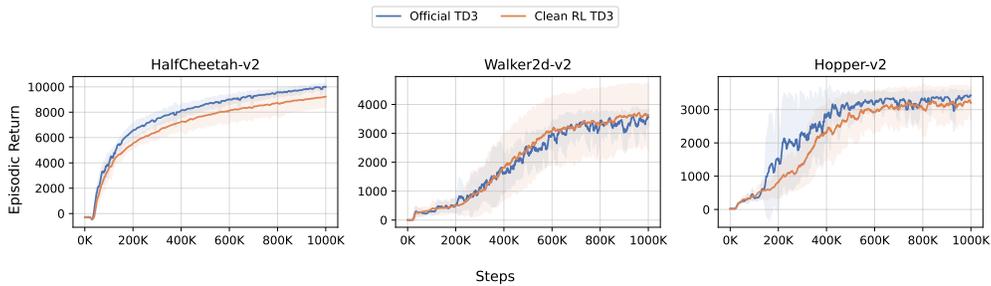


Figure 8: Comparing CleanRL’s TD3 against the original TD3 implementation (sample efficiency).

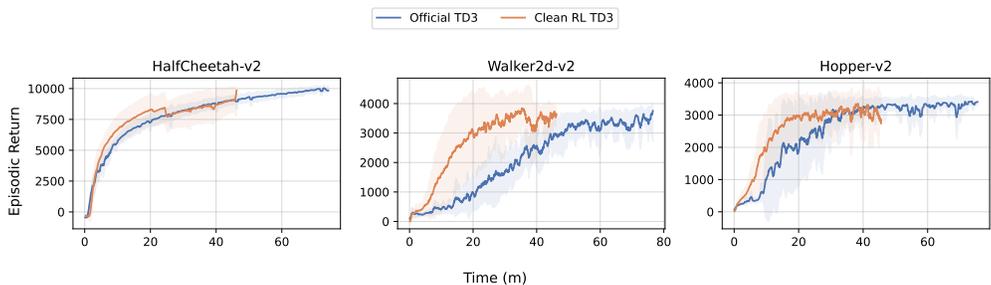


Figure 9: Comparing CleanRL’s TD3 against the original TD3 implementation (time).

684 A.1.2 RLiabLe integration

685 ORLB also integrates with RLiabLe (Agarwal et al., 2021). To enable such plot, the option `--rliable`
 686 can be toggled, then additional parameters are available under `--rc`. Figures 10, 11, 12, 13 showcase
 687 the resulting plots of the following command:

```

688 python -m openrlbenchmark.rlops \
689     --filters '?we=openrlbenchmark&wpm=baselines&ceik=env&cen=exp_name&metric=charts/episodic_return' \
690     baselines-ppo2-cnn?cl=OpenAI Baselines PPO2' \
691     --filters '?we=openrlbenchmark&wpm=envpool-atari&ceik=env_id&cen=exp_name&metric=charts/' \
692     avg_episodic_return' 'ppo_atari_envpool_xla_jax_truncation?cl=CleanRL PPO' \
693     --env-ids AlienNoFrameskip-v4 AmidarNoFrameskip-v4 AssaultNoFrameskip-v4 AsterixNoFrameskip-v4
694     AsteroidsNoFrameskip-v4 AtlantisNoFrameskip-v4 BankHeistNoFrameskip-v4 BattleZoneNoFrameskip-v4
695     BeamRiderNoFrameskip-v4 BerzerkNoFrameskip-v4 BowlingNoFrameskip-v4 BoxingNoFrameskip-v4
696     BreakoutNoFrameskip-v4 CentipedeNoFrameskip-v4 ChopperCommandNoFrameskip-v4
697     CrazyClimberNoFrameskip-v4 DefenderNoFrameskip-v4 DemonAttackNoFrameskip-v4 DoubleDunkNoFrameskip-
698     v4 EnduroNoFrameskip-v4 FishingDerbyNoFrameskip-v4 FreewayNoFrameskip-v4 FrostbiteNoFrameskip-v4
699     GopherNoFrameskip-v4 GravitarNoFrameskip-v4 HeroNoFrameskip-v4 IceHockeyNoFrameskip-v4
700     PrivateEyeNoFrameskip-v4 QbertNoFrameskip-v4 RiverraidNoFrameskip-v4 RoadRunnerNoFrameskip-v4
701     RobotankNoFrameskip-v4 SeaquestNoFrameskip-v4 SkiingNoFrameskip-v4 SolarisNoFrameskip-v4
702     SpaceInvadersNoFrameskip-v4 StarGunnerNoFrameskip-v4 SurroundNoFrameskip-v4 TennisNoFrameskip-v4
703     TimePilotNoFrameskip-v4 TutankhamNoFrameskip-v4 UpNDownNoFrameskip-v4 VentureNoFrameskip-v4
704     VideoPinballNoFrameskip-v4 WizardOfWorNoFrameskip-v4 YarsRevengeNoFrameskip-v4 ZaxxonNoFrameskip-
705     v4 JamesbondNoFrameskip-v4 KangarooNoFrameskip-v4 KrullNoFrameskip-v4 KungFuMasterNoFrameskip-v4
706     MontezumaRevengeNoFrameskip-v4 MsPacmanNoFrameskip-v4 NameThisGameNoFrameskip-v4
707     PhoenixNoFrameskip-v4 PitfallNoFrameskip-v4 PongNoFrameskip-v4 \
708     --env-ids Alien-v5 Amidar-v5 Assault-v5 Asterix-v5 Asteroids-v5 Atlantis-v5 BankHeist-v5 BattleZone-v5
709     BeamRider-v5 Berzerk-v5 Bowling-v5 Boxing-v5 Breakout-v5 Centipede-v5 ChopperCommand-v5
710     CrazyClimber-v5 Defender-v5 DemonAttack-v5 DoubleDunk-v5 Enduro-v5 FishingDerby-v5 Freeway-v5
711     Frostbite-v5 Gopher-v5 Gravitar-v5 Hero-v5 IceHockey-v5 PrivateEye-v5 Qbert-v5 Riverraid-v5
712     RoadRunner-v5 Robotank-v5 Seaquest-v5 Skiing-v5 Solaris-v5 SpaceInvaders-v5 StarGunner-v5
713     Surround-v5 Tennis-v5 TimePilot-v5 Tutankham-v5 UpNDown-v5 Venture-v5 VideoPinball-v5 WizardOfWor-
714     v5 YarsRevenge-v5 Zaxxon-v5 Jamesbond-v5 Kangaroo-v5 Krull-v5 KungFuMaster-v5 MontezumaRevenge-v5
715     MsPacman-v5 NameThisGame-v5 Phoenix-v5 Pitfall-v5 Pong-v5 \
716     --no-check-empty-runs \
717     --pc.ncols 5 \
718     --pc.ncols-legend 2 \
719     --rliable \
720     --rc.score_normalization_method atari \
721     --rc.normalized_score_threshold 8.0 \
722     --rc.sample_efficiency_plots \
723     --rc.sample_efficiency_and_walltime_efficiency_method Median \
724     --rc.performance_profile_plots \
725     --rc.aggregate_metrics_plots \
726     --rc.sample_efficiency_num_bootstrap_reps 50000 \
727     --rc.performance_profile_num_bootstrap_reps 2000 \
728     --rc.interval_estimates_num_bootstrap_reps 2000 \
729     --output-filename static/cleanrl_vs_baselines_atari \
730     --scan-history
  
```

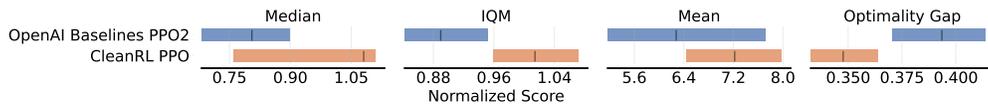


Figure 10: Clean RL PPO vs. OpenAI Baselines PPO, normalized score (RLiabLe).

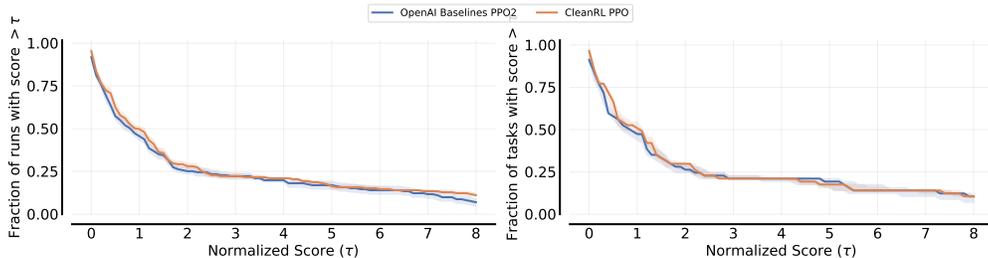


Figure 11: Clean RL PPO vs. OpenAI Baselines PPO, performance profile (RLiabLe).

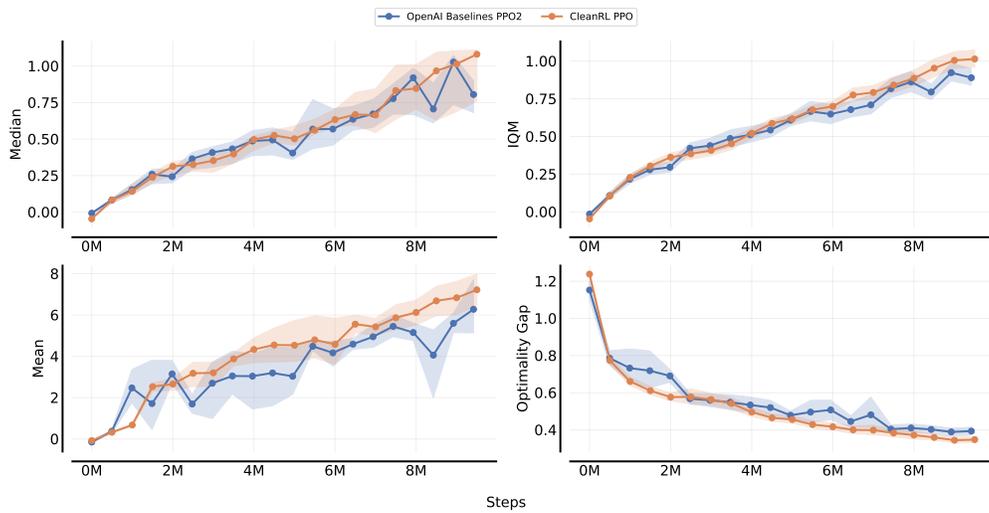


Figure 12: Clean RL PPO vs. OpenAI Baselines PPO, sample efficiency (RLiable).

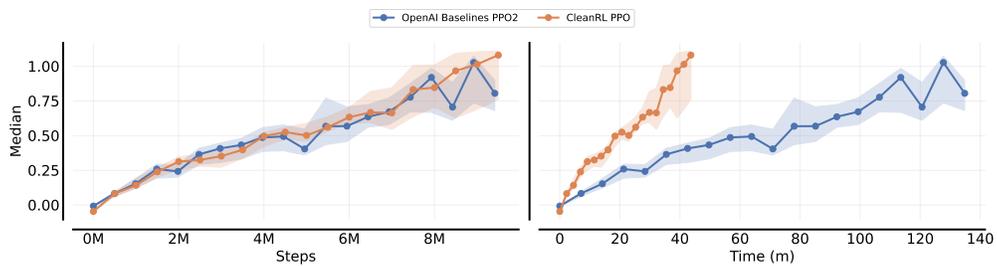


Figure 13: Clean RL PPO vs. OpenAI Baselines PPO, walltime efficiency (RLiable).

731 **A.1.3 Multi-metrics**

732 Sometimes, such as in multi-objective RL (MORL), it is useful to report multiple metrics in the paper.
 733 Hence, the CLI includes an option to plot multiple metrics. Below is an example of CLI and resulting
 734 plots (Figure 14) for multiple MORL algorithms on different environments.

```

735 python -m openrlbenchmark.rlops_multi_metrics \
736 --filters '?we=openrlbenchmark&wpm=MORL-Baselines&ceik=env_id&cen=algo&metrics=eval/hypervolume&metrics=
737 eval/igd&metrics=eval/sparsity&metrics=eval/mul' \
738 'Pareto Q-Learning?cl=Pareto Q-Learning' \
739 'MultiPolicy MO Q-Learning?cl=MPMOQL' \
740 'MultiPolicy MO Q-Learning (OLS)?cl=MPMOQL (OLS)' \
741 'MultiPolicy MO Q-Learning (GPI-LS)?cl=MPMOQL (GPI-LS)' \
742 --env-ids deep-sea-treasure-v0 deep-sea-treasure-concave-v0 fruit-tree-v0 \
743 --pc.ncols 3 \
744 --pc.ncols-legend 4 \
745 --pc.xlabel 'Training steps' \
746 --pc.ylabel '' \
747 --pc.max_steps 400000 \
748 --output-filename morl/morl_deterministic_envs \
749 --scan-history
    
```

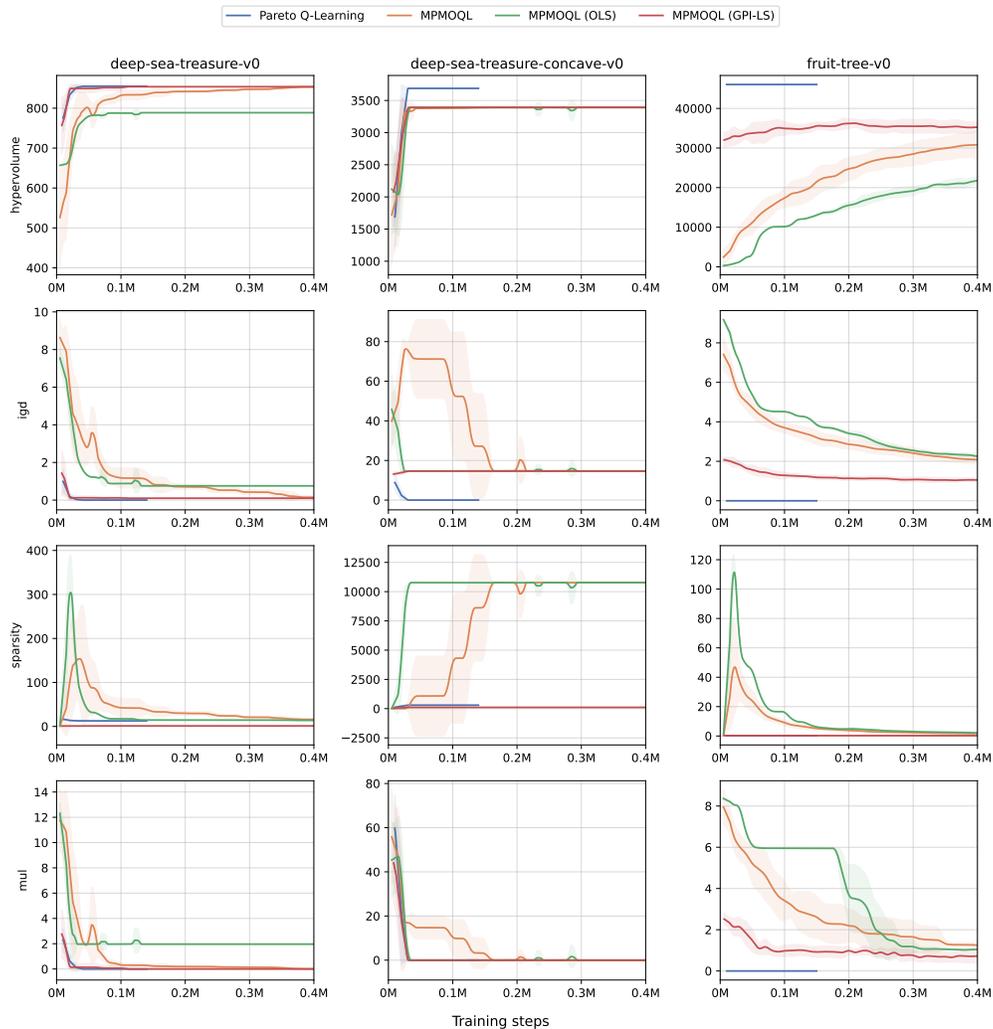


Figure 14: Plotting different metrics for different environments.

750 A.2 Using a custom script

751 Our CLI proves highly beneficial for generating standard RL plots, as demonstrated above. Neverthe-
752 less, in certain specialized cases, researchers may wish to expose the data in an alternative format.
753 Fortunately, all the data hosted in ORLB is accessible through the Weights and Biases API. The
754 following example illustrates how this API can be utilized. From there, researchers can employ any
755 custom script for plotting this data to suit their specific needs. A simple example of such a script is
756 given below, and the corresponding generated plot is shown in Figure 15.

```
757 import matplotlib.pyplot as plt
758 import wandb
759
760 project_name = "sb3"
761 run_id = "0a1kqgev"
762
763 api = wandb.Api()
764 run = api.run(f"openrlbenchmark/{project_name}/{run_id}")
765 history = run.history(keys=["global_step", "eval/mean_reward"])
766 plt.plot(history["global_step"], history["eval/mean_reward"])
767 plt.title(run.name)
768 plt.savefig("custom_plot.png")
```

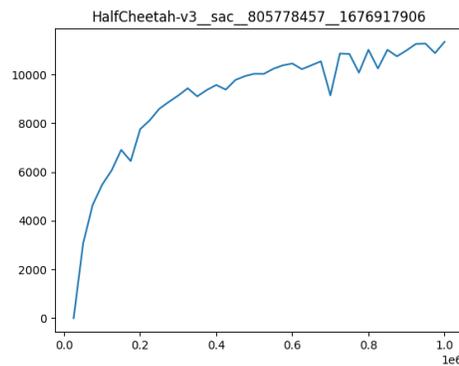


Figure 15: Example of a plot created with a custom script, by importing data directly from ORLB using the WandB API.

769 B Additional details for the case study

770 This appendix gives additional results related to the first case study presented in Section 3.1. Figure
771 17 shows the results by environment for the Atari benchmark, and Figure 16 shows them for the
772 MuJoCo and Box2d benchmarks. The command lines used to generate these figures are as follows.

```
773 python -m openrlbenchmark.rlops \  
774 --filters '?we=openrlbenchmark&wpn=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'ppo?cl=PP0' \  
775 --filters '?we=modanesh&wpn=openrlbenchmark&ceik=env&cen=algo&metric=eval/mean_reward' 'ppo?cl=PP0 w/  
776 MC for value estimation' \  
777 --env-ids BreakoutNoFrameskip-v4 SpaceInvadersNoFrameskip-v4 SeaquestNoFrameskip-v4 EnduroNoFrameskip-  
778 v4 PongNoFrameskip-v4 QbertNoFrameskip-v4 BeamRiderNoFrameskip-v4 \  
779 --no-check-empty-runs \  
780 --pc.ncols 3 \  
781 --pc.ncols-legend 2 \  
782 --rliable \  
783 --rc.score_normalization_method atari \  
784 --rc.normalized_score_threshold 8.0 \  
785 --rc.sample_efficiency_plots \  
786 --rc.sample_efficiency_and_walltime_efficiency_method Median \  
787 --rc.performance_profile_plots \  
788 --rc.aggregate_metrics_plots \  
789 --rc.sample_efficiency_num_bootstrap_reps 1000 \  
790 --rc.performance_profile_num_bootstrap_reps 1000 \  
791 --rc.interval_estimates_num_bootstrap_reps 1000 \  
792 --output-filename static/gae_for_ppo_value_atari_per_env \  
793 --scan-history \  
794 --rc.sample_efficiency_figsize 7 4  
795  
796 python -m openrlbenchmark.rlops \  
797 --filters '?we=openrlbenchmark&wpn=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'ppo?cl=PP0' \  
798 --filters '?we=modanesh&wpn=openrlbenchmark&ceik=env&cen=algo&metric=eval/mean_reward' 'ppo?cl=PP0 w/  
799 MC for value estimation' \  
800 --env-ids InvertedDoublePendulum-v2 InvertedPendulum-v2 Reacher-v2 HalfCheetah-v3 Hopper-v3 Swimmer-v3  
801 Walker2d-v3 LunarLander-v2 \  
802 --no-check-empty-runs \  
803 --pc.ncols 3 \  
804 --pc.ncols-legend 2 \  
805 --rliable \  
806 --rc.normalized_score_threshold 1.0 \  
807 --rc.sample_efficiency_plots \  
808 --rc.sample_efficiency_and_walltime_efficiency_method Median \  
809 --rc.performance_profile_plots \  
810 --rc.aggregate_metrics_plots \  
811 --rc.sample_efficiency_num_bootstrap_reps 1000 \  
812 --rc.performance_profile_num_bootstrap_reps 1000 \  
813 --rc.interval_estimates_num_bootstrap_reps 1000 \  
814 --output-filename static/gae_for_ppo_value_mujoco_per_env \  
815 --scan-history \  
816 --rc.sample_efficiency_figsize 7 4
```

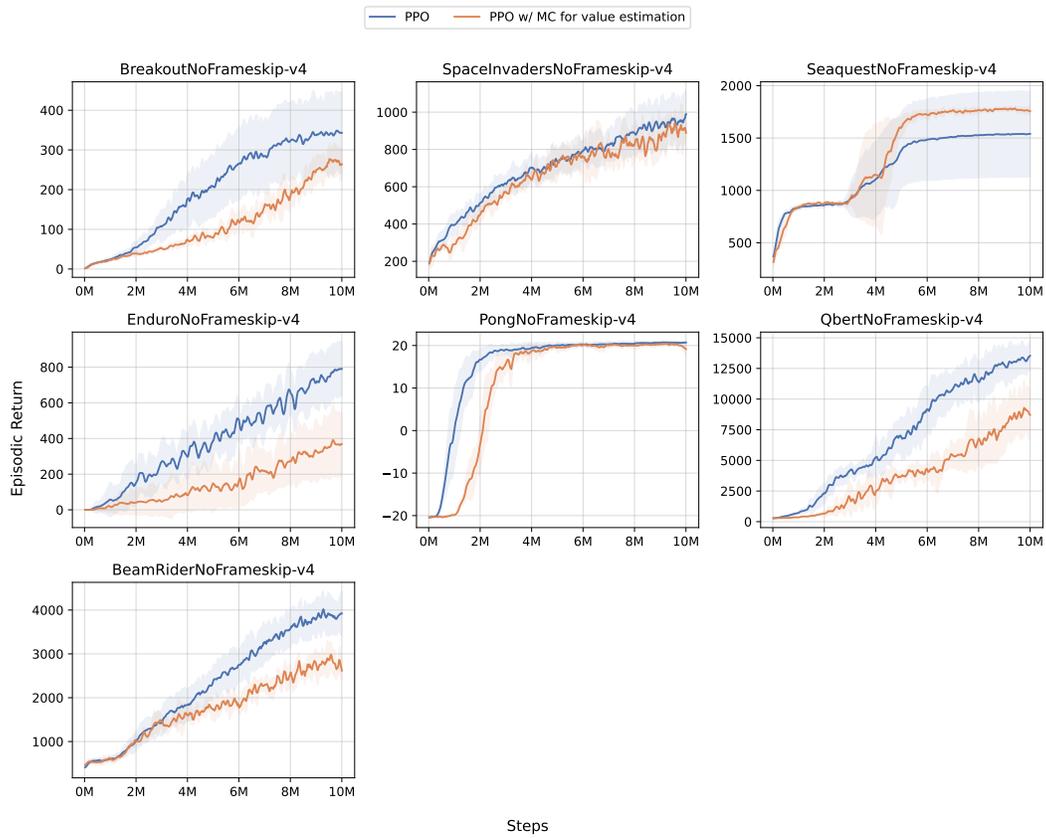


Figure 16: Comparison between the original PPO and the PPO with MC value estimates in various MuJoCo and Box2D environments. Plots represent the evolution of the episodic return as a function of the number of interactions with the environment, and shaded areas represent the standard deviation.

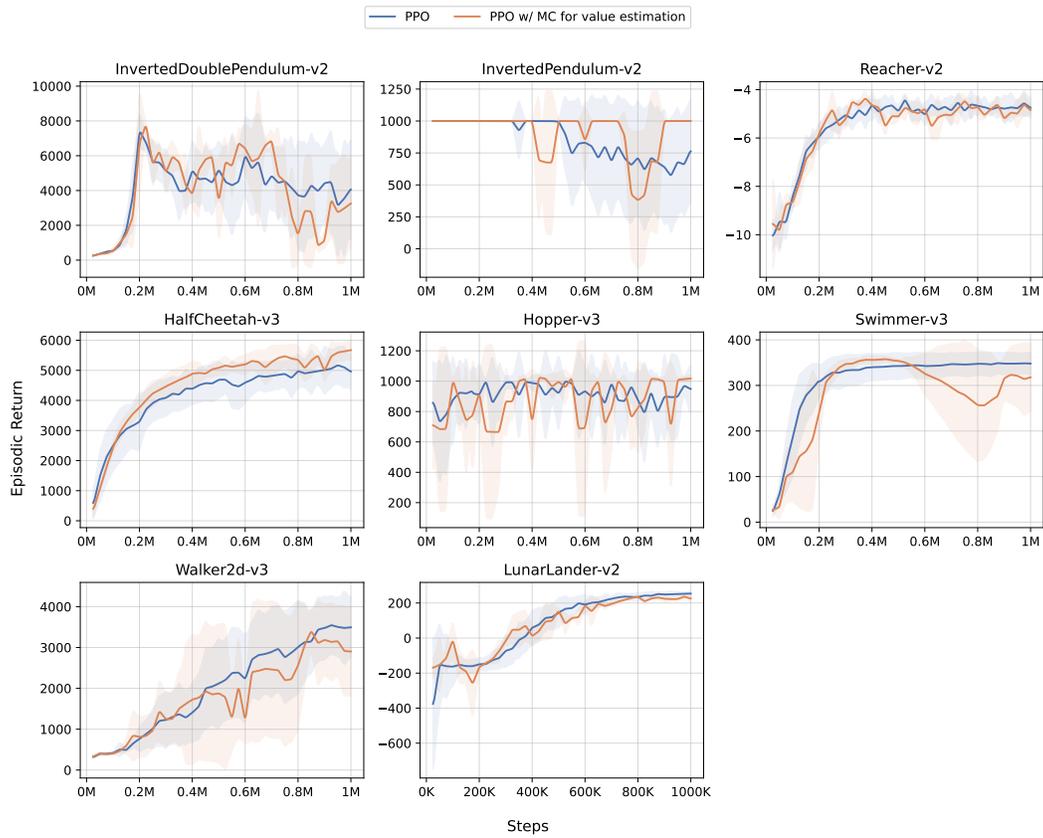


Figure 17: Comparison between the original PPO and the PPO with MC value estimates in various MuJoCo and Box2D environments. Plots represent the evolution of the episodic return as a function of the number of interactions with the environment, and shaded areas represent the standard deviation.

817 C Refine the MuJoCo benchmark with Stable Baselines3

818 In this appendix, we present a synthetic representation of the learning results of the Stable Baselines3
 819 algorithms (Raffin et al., 2021) tested on the MuJoCo benchmark (Brockman et al., 2016; Todorov
 820 et al., 2012), whose data is contained in ORLB. At the time of writing, data from 757 runs has
 821 been used, unevenly distributed between the different experiments. It is important to emphasise that
 822 the optimisation of hyperparameters and the training budget vary from one experiment to another.
 823 Consequently, the results should be interpreted with caution. All the hyperparameters and raw
 824 data used to generate these curves are available on ORLB. Figure 18 shows the aggregation of the
 825 final performances following the recommendations of Agarwal et al. (2021), and Figure 19 the
 826 corresponding performance profiles. Figure 20 shows the learning curves as a function of the number
 827 of interactions.

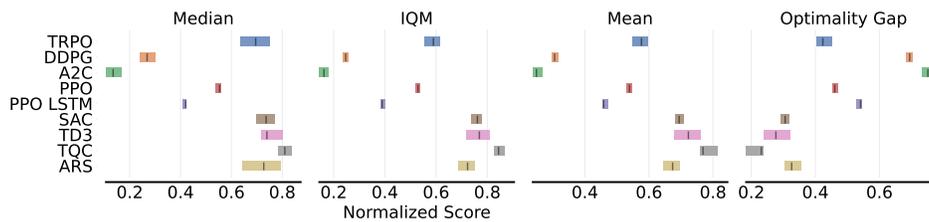


Figure 18: Aggregated final normalized episodic return with 95% stratified bootstrap CIs on the MuJoCo benchmark of the algorithms integrated into Stable Baselines3.

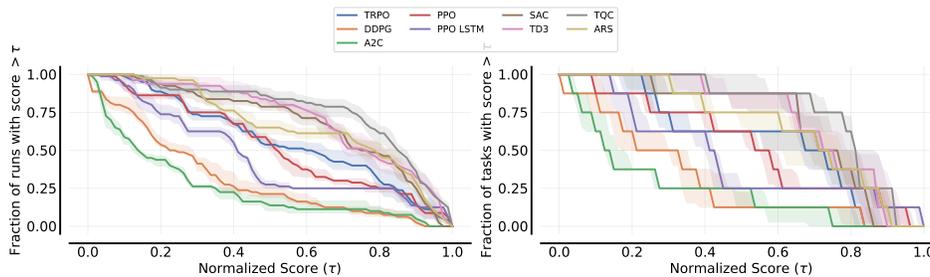


Figure 19: Performance profile of algorithms implemented using Stable Baselines 3 (Raffin et al., 2021) on the MuJoCo benchmark (Todorov et al., 2012). Scores are normalized using the min-max method.

828 The command used to generate Figures 18, 19 and 20 is as follows⁷.

```
829 python -m openrlbenchmark.rlops \
830 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'trpo?cl=TRPO' \
831 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'ddpg?cl=DDPG' \
832 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'a2c?cl=A2C' \
833 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'ppo?cl=PPO' \
834 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'ppo_lstm?cl=PPO LSTM'
835 \
836 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'sac?cl=SAC' \
837 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'td3?cl=TD3' \
838 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'ars?cl=ARS' \
839 --filters '?we=openrlbenchmark&wps=sb3&ceik=env&cen=algo&metric=eval/mean_reward' 'tqc?cl=TQC' \
840 --env-ids Ant-v3 BipedalWalker-v3 BipedalWalkerHardcore-v3 HalfCheetah-v3 Hopper-v3 Humanoid-v3 Swimmer
841 -v3 Walker2d-v3 \
842 --no-check-empty-runs \
843 --pc.ncols 2 \
844 --pc.ncols-legend 4 \
845 --rliable \
846 --rc.normalized_score_threshold 1.0 \
847 --output-filename static/mujoco_sb3 \
848 --scan-history
```

⁷For Figure 20, we are omitting ARS as it was run with many more steps, and its inclusions hinder readability.

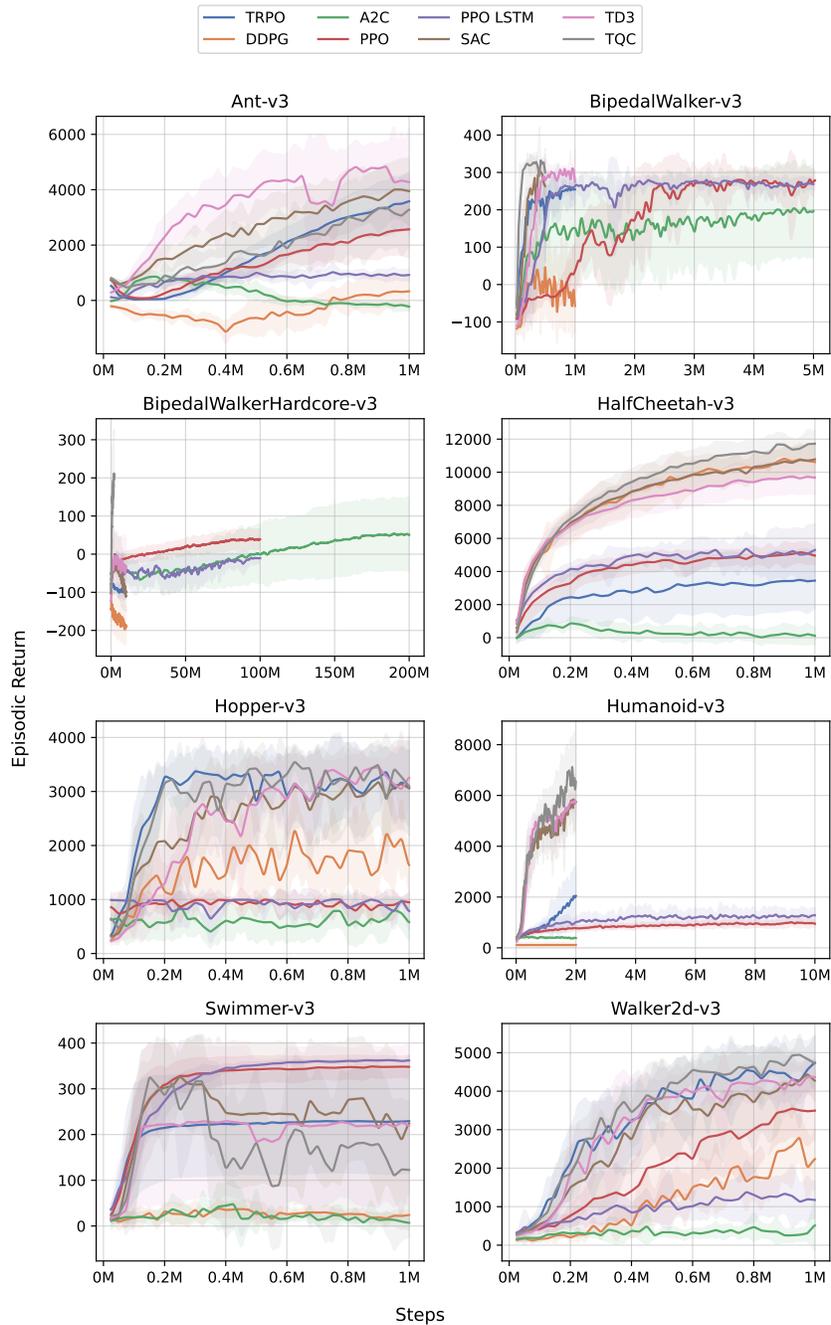


Figure 20: Sample efficiency curves for algorithms on the MuJoCo Benchmark (Todorov et al., 2012). This graph presents the mean episodic return for algorithms implemented using Stable Baselines 3 (Raffin et al., 2021), averaged across a minimum of 10 runs (refer to ORLB for specific run counts). Data points are subsampled to 10,000 and interpolated for clarity. The curves are smoothed using a rolling average with a window size of 100. The shaded regions around each curve indicate the standard deviation.