

CleanRL: High-quality Single-file Implementations of Deep Reinforcement Learning Algorithms

Shengyi Huang¹

Rousslan Fernand Julien Dossa²

Chang Ye³

Jeff Braga¹

COSTA.HUANG@OUTLOOK.COM

DOSS@AI.CS.KOBE-U.AC.JP

C.YE@NYU.EDU

JEFFREYBRAGA@GMAIL.COM

¹College of Computing and Informatics, Drexel University, USA

²Graduate School of System Informatics, Kobe University, Japan

³Tandon School of Engineering, New York University, USA

Editor: PLACEHOLDER

Abstract

CleanRL is an open-source library that provides high-quality single-file implementations of Deep Reinforcement Learning algorithms. It provides a simpler yet scalable developing experience by having a straightforward codebase and integrating production tools to help interact and scale experiments. In **CleanRL**, we put all details of an algorithm into a single file, making these performance-relevant details easier to recognize. Additionally, an experiment tracking feature is available to help log metrics, hyperparameters, videos of an agent’s gameplay, dependencies, and more to the cloud. Despite succinct implementations, we have also designed tools to help scale, at one point orchestrating experiments on more than 2000 machines simultaneously via Docker and cloud providers. Finally, we have ensured the quality of the implementations by benchmarking against a variety of environments.

1. Introduction

In recent years, Deep Reinforcement Learning (DRL) algorithms have achieved great success in training autonomous agents for tasks ranging from playing video games directly from pixels to robotic control (Mnih et al., 2013; Lillicrap et al., 2015; Schulman et al., 2017). At the same time, open-source DRL libraries also flourish in the community, such as **Stable Baselines 3 (SB3)** (Raffin et al., 2019), **RLlib** (Liang et al., 2018), **MushroomRL** (D’Eramo et al., 2020), **PFRL** (Fujita et al., 2021), and others. Many of them have adopted good modular designs and fostered vibrant development communities. Nevertheless, understanding all implementation details of an algorithm remains difficult because these details are spread to different modules. However, it has become increasingly important to understand these details because implementations matter (Engstrom et al., 2019).

In this paper, we introduce **CleanRL**, a DRL library based on single-file implementations to help researchers understand all algorithms’ details, prototype new features, analyze experiments, and scale with ease. **CleanRL** is a *non-modular* library. Each algorithm in **CleanRL** is self-contained in a single file. We have trimmed the lines of code (LOC) to a minimal degree while maintaining readability and decent performance w.r.t reputable reference implementations. For instance, our Proximal Policy Optimization (PPO) (Schulman

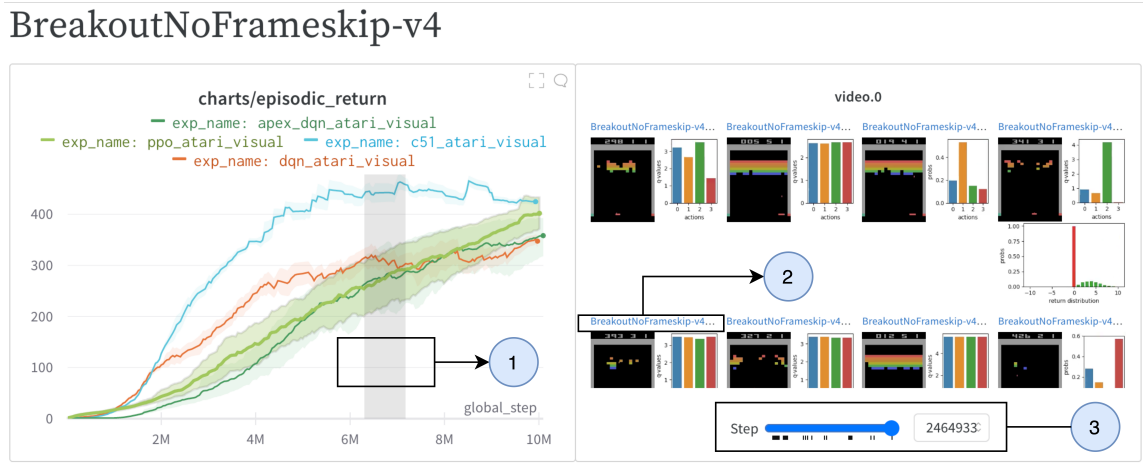


Figure 1: A screenshot of the Open RL Benchmark where the users can (1) zoom into the learning curve, (2) click to trace back to the original experiment for more info (e.g. hyperparameters), and (3) checkout videos at different training stages.

et al., 2017) implementation matches SB3’s PPO performance in the Atari game Breakout using only 337 LOC in a single file¹, making it much easier to understand the algorithm in one go. In contrast, the researchers of other libraries need to understand the modularity design (on average 7-15 files) which can be comprised of thousands of LOC.

While CleanRL is stand-alone, for convenience, it supports optional integration with production-quality tool providers. Firstly, integrating **Weights and Biases (W&B)** helps track the logs and metrics to the cloud. Over the years, we have tracked thousands of experiments across 7+ algorithms and 40+ games in our Open RL Benchmark (<https://benchmark.cleanrl.dev>) to ensure and showcase the quality of our implementations (see Figure 1). Secondly, integrating **AWS Batch** helps provision any AWS instances as computational nodes for large-scale experiment orchestration. Adopting these tools has helped CleanRL focus just on algorithms without worrying about analysis or scaling utilities. Also, our architecture design allows the users to flexibly use other tracking and orchestration vendors if desired, so there is no dependency on proprietary services.

2. Single-file Implementations

Even small implementation details can hugely impact the performance of deep RL algorithms (Engstrom et al., 2019), however the nature of the modular libraries encapsulates and hides these details. CleanRL helps the researchers to recognize these details by adopting single-file implementations. That is, we put the **algorithm details** and the **environment-specific details** into a single file. For example, we have a

1. `ppo.py` (321 LOC) for the classic control environments, such as `CartPole-v1`.

1. https://github.com/vwxyzjn/cleanrl/blob/2486b90102bd87fa43cda481fe711e0eda0c162b/cleanrl/ppo_atari.py

2. `ppo_atari.py` (337 LOC) for the Atari environments (Bellemare et al., 2013).
3. `ppo_continuous_action.py` (331 LOC) for the robotics environments (e.g. MuJoCo, PyBullet) with continuous action spaces (Schulman et al., 2017).

Despite having duplicate code among these files, the single-file implementations have the following benefits.

1. **Transparent learning experience** It becomes easier to recognize all aspects of the code in one place. By looking at `ppo.py`, it is straightforward to recognize the 11 general implementation details of PPO, especially when watching the corresponding tutorial². Then, comparing it with `ppo_atari.py` shows about 30 LOC difference to implement 9 Atari-specific details³, and another comparison with `ppo_continuous_action.py` shows about 25 LOC difference to implement 6 details for continuous action spaces⁴.
2. **Faster debug experience** Most variables in the files live in the *global python name scope*. This means the researchers can do `Ctrl+C` to stop the program execution and check most variables and their shapes in the interactive shell (see Appendix A). This is arguably more convenient than using the Python’s debugger, which only shows the variables in a specific name scope like in a function.
3. **Faster prototype experience** Because of the faster debug experience, it becomes easier to develop new features. For example, invalid action masking (Huang and Ontañón, 2020) is a common technique used in games with large parameterized action spaces. With `CleanRL`, it takes about 40 LOC to implement⁵. In comparison, it takes substantially more LOC (e.g. 600+ LOC not counting test cases) because of overhead such as re-factoring the functional arguments, making the classes more general.

For this motivation, we have also implemented `dqn.py`, `dqn_atari.py` (Mnih et al., 2013), `c51.py`, `c51_atari.py` (Bellemare et al., 2017), `apex_atari.py` (Horgan et al., 2018), `ddpg_continuous_action.py` (Lillicrap et al., 2015), `td3_continuous_action.py` (Fujimoto et al., 2018), and `sac_continuous_action.py` (Haarnoja et al., 2018).

3. Experiment Tracking and Analysis

`CleanRL` supports integration with `W&B`, a modern experiment tracking and analysis solution. When running an experiment, `W&B` automatically tracks the 1) source code, 2) dependency (i.e. `requirements.txt`), 3) hyperparameters, 4) the exact terminal command used to run the experiment, 5) training metrics, 6) videos of the agents playing the game, 7) system metrics, and 8) logs (`stdout`, `stderr`). See Appendix B for 9 selected features that help query and clean the experiment data.

Although `W&B` is a proprietary service that we elect to use, our implementation uses `Tensorboard` to record metrics, which is open-source. So the users could easily change about 10 LOC to use other experiment tracking vendors such as `ClearML` (<https://clear.ml>) that supports `Tensorboard` integration.

2. See <https://youtu.be/MEt6rrxH8W4>

3. See https://youtu.be/05RMTj-2K_Y

4. See https://youtu.be/_rDUE0sjUds

5. <https://www.diffchecker.com/g9TYkKCS>

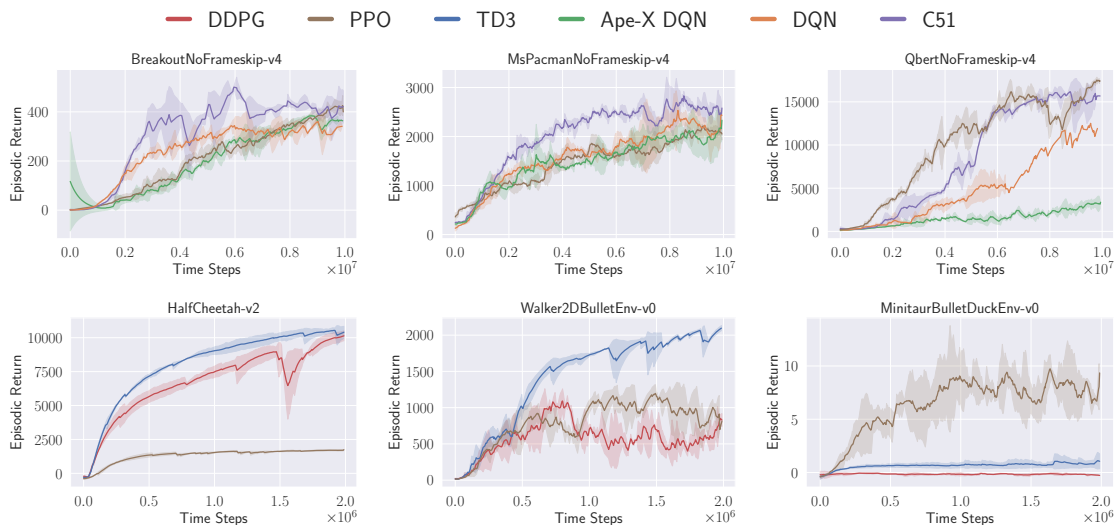


Figure 2: The performance of `CleanRL`'s algorithms in Atari and continuous action domains.

4. Cloud Integration

Despite succinct implementations, we are able to use `CleanRL` for large-scale studies by leveraging public cloud infrastructure such as Amazon Web Services (AWS). We package the code into a `Docker` container to orchestrate the experiments on any cloud provider. When needing to run experiments on demand, the users of `CleanRL` can use `Terraform` to reproduce a preset infrastructure in a few commands, then utilize our utility scripts to create as many experiments as needed. In 2020 alone, the authors have run over 50,000+ hours of experiments using this workflow. See <https://docs.cleanrl.dev> for more documentation.

5. Open RL Benchmark

We have created a project called Open RL Benchmark (<https://benchmark.cleanrl.dev>) to track canonical RL experiments. To date, we have tracked thousands of experiments across 7+ algorithms in various domains (e.g. Atari, MuJoCo, PyBullet, Procgen (Cobbe et al., 2020), Griddy (Bamford et al., 2020), Gym- μ RTS (Huang et al., 2021)). See Figure 2 for selected experiments. Our benchmark is interactive as shown in Figure 1: researchers can easily query information such as GPU utilization and videos of an agent's gameplay that are normally hard to acquire in other RL benchmarks.

6. Conclusion

In this paper, we introduced `CleanRL` as a deep RL library that provides high-quality single-file implementations with experiment tracking and cloud integration. The source code can be found at <https://github.com/vwxyzjn/cleanrl>.

References

- Chris Bamford, Shengyi Huang, and Simon Lucas. Griddly: A platform for ai research in games, 2020.
- M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47: 253–279, jun 2013.
- Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning*, pages 449–458. PMLR, 2017.
- Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In *International conference on machine learning*, pages 2048–2056. PMLR, 2020.
- Carlo D’Eramo, Davide Tateo, Andrea Bonarini, Marcello Restelli, and Jan Peters. Mushroomrl: Simplifying reinforcement learning research. *Journal of Machine Learning Research*, 2020.
- Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. Implementation matters in deep rl: A case study on ppo and trpo. In *International Conference on Learning Representations*, 2019.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596. PMLR, 2018.
- Yasuhiro Fujita, Prabhat Nagarajan, Toshiki Kataoka, and Takahiro Ishikawa. Chainerrl: A deep reinforcement learning library. *Journal of Machine Learning Research*, 22(77): 1–14, 2021.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018.
- Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado Van Hasselt, and David Silver. Distributed prioritized experience replay. *arXiv preprint arXiv:1803.00933*, 2018.
- Shengyi Huang and Santiago Ontañón. A closer look at invalid action masking in policy gradient algorithms. *arXiv preprint arXiv:2006.14171*, 2020.
- Shengyi Huang, Santiago Ontañón, Chris Bamford, and Lukasz Grela. Gym- μ rts: Toward affordable full game real-time strategy games research with deep reinforcement learning. *arXiv preprint arXiv:2105.13807*, 2021.

Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph Gonzalez, Michael Jordan, and Ion Stoica. Rllib: Abstractions for distributed reinforcement learning. In *International Conference on Machine Learning*, pages 3053–3062. PMLR, 2018.

Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous Control with Deep Reinforcement Learning. *arXiv preprint arXiv:1509.02971*, 2015.

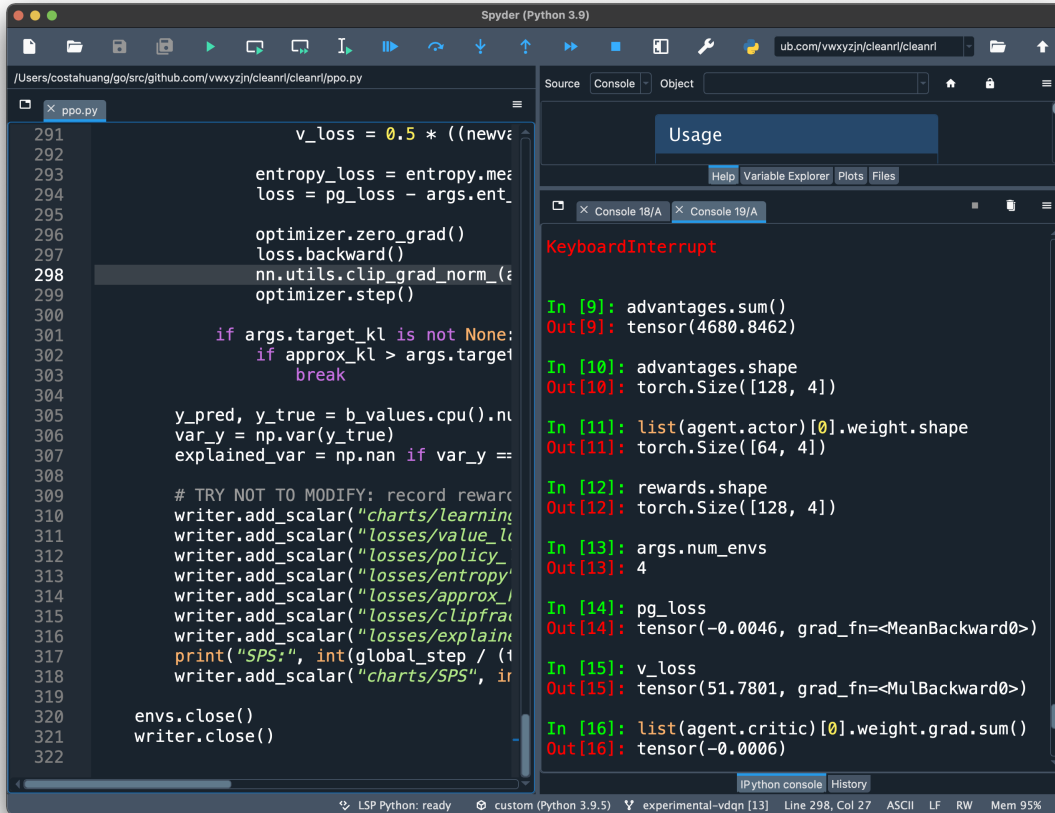
Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing Atari with Deep Reinforcement Learning. *arXiv preprint arXiv:1312.5602*, 2013.

Antonin Raffin, Ashley Hill, Maximilian Ernestus, Adam Gleave, Anssi Kanervisto, and Noah Dormann. Stable baselines3. <https://github.com/DLR-RM/stable-baselines3>, 2019.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Appendix A. Interactive Shell

In CleanRL, we have put most of the variables in the *global python name scope*. This makes it easier to inspect the variables and their shapes. The following figure shows a screenshot of the Spyder editor ⁶, where the code is on the left and the interactive shell is on the right. In the interactive shell, we can easily inspect the variables for debugging purposes without modifying the code.



6. <https://www.spyder-ide.org/>

Appendix B. W&B Editing Panel

A screenshot of the W&B panel that allows the users to change smoothing weight, add panels to show different metrics like losses, visualize the videos of the agents' gameplay, filter, group, sort, and search for desired experiments.

