
PufferLib: Making Reinforcement Learning Libraries and Environments Play Nice

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Common simplifying assumptions often cause standard reinforcement learning
2 (RL) methods to fail on complex, open-ended environments. Creating a new wrap-
3 per for each environment and learning library can help alleviate these limitations,
4 but building them is labor-intensive and error-prone. This practical tooling gap
5 restricts the applicability of RL as a whole. To address this challenge, Puffer-
6 Lib transforms complex environments into a broadly compatible, vectorized for-
7 mat that eliminates the need for bespoke conversion layers and enables rigorous
8 cross-environment testing. PufferLib does this without deviating from standard
9 reinforcement learning APIs, significantly reducing the technical overhead. We
10 release¹ PufferLib’s complete source code under the MIT license, a pip module, a
11 containerized setup, comprehensive documentation, and example integrations. We
12 also maintain a community Discord channel to facilitate support and discussion.

13 1 Background and Introduction

14 Reinforcement Learning (RL) generates data through interaction with a multitude of parallel en-
15 vironment simulations. This dynamism introduces non-stationarity into the optimization process,
16 necessitating algorithmic treatments distinct from those employed in supervised learning. When
17 compounded by sparse reward signals, this issue yields several complications, including extreme
18 sensitivity to hyperparameters, which extends even to the random seed. Consequently, experiments
19 often yield unpredictable learning curves with spikes, plateaus, or crashes, deviating from the more
20 reliable behavior observed in other machine learning domains.

21 Alongside this lies the pragmatic challenge of implementing RL in complex environments with
22 currently available tools. Although this is arguably a more solvable problem than optimizing the
23 online learning process, the lack of effective tooling often exacerbates the problem, making it an
24 arduous task to resolve despite thorough investigation. These issues frequently cause significant
25 delays, frustration, and stagnation in the field, potentially deterring talented researchers from pursuing
26 work in this area.

27 PufferLib is a middleware layer that resolves longstanding compatibility issues between environments
28 and tools. **Our work enables other research in open-endedness by removing a large practical
29 barrier to studying complex environments.** Existing solutions such as Gym/Gymnasium [Brockman
30 et al., 2016], PettingZoo [Terry et al., 2020b], and SuperSuit [Terry et al., 2020a] aim to define standard

¹As this is a software release with an existing community, it is impossible to fully anonymize the submission. We ask reviewers to take this practical limitation into account when assessing our adherence to double blind.

The PufferLib System Architecture for Broad Environment Compatibility

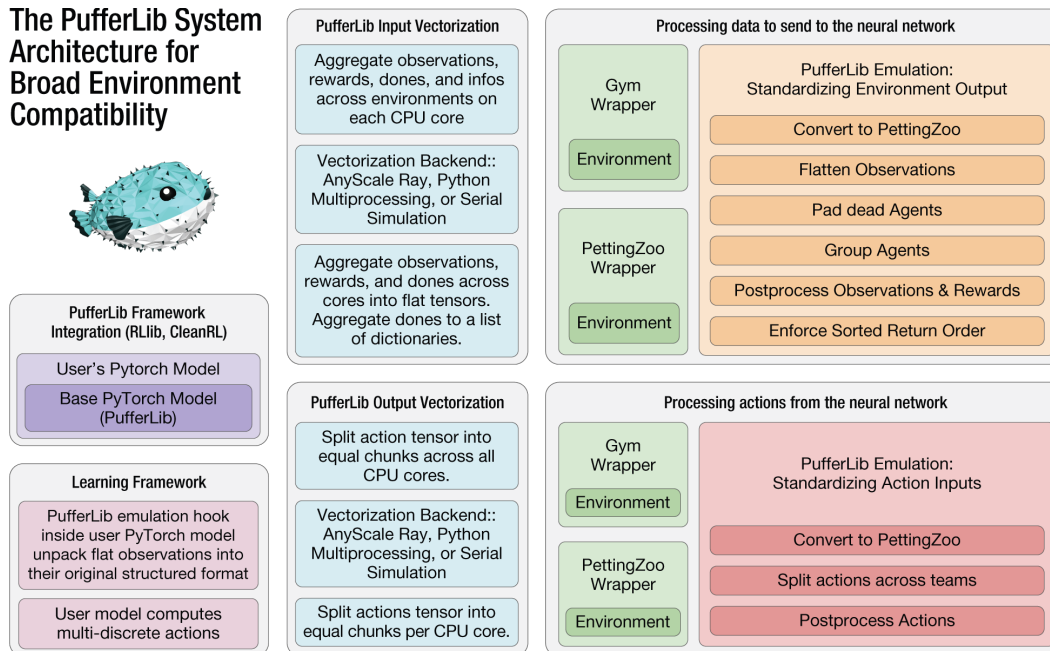


Figure 1: Detailed but non-comprehensive illustration of the PufferLib system architecture, comprising emulation, vectorization, and learning framework integrations. The orange emulation block demonstrate how PufferLib receives and processes environment data. The red emulation block demonstrates how PufferLib processes actions from the neural network to send to the environment. The blue vectorization blocks aggregate and split data received from and sent to the environment. Finally, the pink and purple blocks summarize how PufferLib provides compatibility with multiple frameworks given a single PyTorch network.

31 APIs for environments and implement common wrappers. PufferLib builds on Gym and PettingZoo
 32 but also addresses their specific limitations, which we will discuss after providing comprehensive
 33 context for the problem at hand.

34 PufferLib allows users to wrap most new environments in a single line of code, without changing
 35 the original Gym/Gymnasium/PettingZoo API. This wrapper makes environments compatible with
 36 popular reinforcement learning libraries, such as CleanRL [Huang et al., 2021a] and RLLib [Liang
 37 et al., 2017]. It natively supports multi-agent and variable-agent environments and addresses common
 38 complexities that include batching structured observations and actions, featurizing observations,
 39 shaping rewards, and grouping agents into teams. PufferLib is also designed for extensibility and is
 40 capable of supporting new learning libraries with a complete feature set in typically about a hundred
 41 lines of code.

42 2 Problem Statement

43 Middleware limitations are seldom topics of sustained academic interest, but perhaps for this reason,
 44 their impact can be prolonged and severe. To thoroughly ground our work, we will walk through the
 45 intricacies of the transformations that reinforcement learning data must undergo, and demonstrate the
 46 shortcomings of existing approaches. Specifically, we will trace the required transformations from
 47 simulation onset to data processing by the initial model layer, and from action computation to the
 48 point when those actions influence the environment.

49 We will use Neural MMO [Suarez et al., 2021], a PettingZoo-compliant environment, as our guiding
 50 example. This environment encapsules many complexities common in advanced environments. It
 51 features 128 agents competing to complete open-ended tasks in a procedurally-generated world.

52 Agents are provided with rich, structured observations of their surroundings and a hierarchical action
53 space for interactions.

54 The environment initialization starts with a configuration file and a reset to yield an initial set of
55 observations. This results in a dictionary of 128 individual observations, each of which is a structured
56 dictionary housing differently-shaped tensors related to various aspects of the observation. As a part
57 of the environment’s standard training setting, these agents are grouped into teams of 8. Each team
58 observation is then processed by a featurizer to yield a single structured observation, aggregating
59 information from across the team’s agents. Subsequently, this observation must be batched for model
60 usage.

61 This introduces two challenges. First, since the observation is structured, we cannot merely con-
62 catenate tensors; we must concatenate each sub-observation across agents. Secondly, many learning
63 libraries presuppose that observations can be stored in flat tensors, thus requiring data flattening.
64 Following this, the data must be concatenated with information from several parallel environment
65 instances. Once done, the data can be forwarded to the network.

66 We now encounter another problem: the network itself is structured, and attempting to learn from
67 the flattened representation is akin to unraveling an image and using dense layers. Therefore, the
68 structured observation representation must be recovered in a batched form, allowing for efficient
69 processing of each sub-observation across all teams and environments in parallel. The model
70 then computes a multidiscrete output distribution and samples an integer array for each team and
71 environment. The output data is divided across environments, and each multidiscrete action is mapped
72 into a structured format where each integer signifies a specific agent’s action within a team. Finally,
73 the environment can execute its first step.

74 Regrettably, this is the least complex step. All preceding actions must be reiterated, but with additional
75 complexities. For example, the environment must now also return *rewards*, *done*s, and *infos*. These
76 outputs, particularly rewards and *done*s, require grouping by team. For each team, we must track
77 which agents have completed their tasks and signal that team is *done* only when all agents have
78 finished. Similarly, we need a method to post-process and group *reward* signals per team. Since most
79 learning libraries anticipate each agent to return an observation on every step, we must zero-pad the
80 tensor for any agents that are *done*. Moreover, as the PettingZoo API does not mandate a consistent
81 observation return order (a common source of bugs), we must verify this as well.

82 As illustrated, considerable work is needed to ensure compatibility between the environment and
83 standard learning libraries - even for a fully Gym and PettingZoo-compliant environment like Neural
84 MMO. We have provided support to the Neural MMO team in integrating PufferLib, and prior to
85 integration, about a quarter of the Neural MMO code base was devoted to these transformations. This
86 was also the primary source of bugs, many of which would lead to silent performance degradations.
87 For instance, specific patterns of agent deaths could cause incorrectly ordered observations, leading
88 to neural networks optimizing trajectories assembled from different agents. In another case, a bug in
89 the reconstruction of observations misaligned data, causing incorrect subnetwork processing. Despite
90 a strong engineering focus on testing, these bugs are two among dozens that reportedly emerged
91 during Neural MMO’s development.

92 **3 Related Tools**

93 Gym and PettingZoo, the prevalent environment APIs for single-agent and multi-agent environments
94 respectively, offer several tools to mitigate the complexities described earlier. Supplementary third-
95 party tools, like SuperSuit, provide standalone wrappers, while numerous reinforcement learning
96 libraries furnish wrappers compatible with their internal APIs. For instance, Gym provides a
97 range of wrappers for image observation preprocessing, observation flattening, action and reward
98 postprocessing, and even sanity check wrappers for bug prevention. SuperSuit further adds multi-
99 agent wrappers specifically designed to address the agent termination and padding issues discussed
100 previously.

101 Current methodologies present some significant challenges. The tools described are designed as a set
102 of wrappers applied sequentially to an environment instance, implying that (with a few exceptions),
103 they should function in any order. However, particularly with PettingZoo, which caters to multi-
104 and variable-agent environments, the gamut of possible environments is vast and challenging to test.
105 This often results in scenarios where a bug in one wrapper causes an error in a different wrapper.
106 Identifying the origin of such errors across deeply nested wrapper classes can be an overwhelming
107 task, contributing to a general sense of frustration common in reinforcement learning research.

108 Moreover, the coverage of wrappers is insufficient. Despite the difficulties in testing and maintaining
109 compatibility among existing wrappers, more are still needed. As it stands, there is no wrapper
110 ensuring consistent agent key ordering, despite many reinforcement learning libraries demanding this.
111 No wrapper exists for grouping agents into teams, a common operation, nor a wrapper that natively
112 vectorizes multi-agent environments across multiple cores. The current workarounds for the latter are
113 unstable, abusing single-agent vectorization code. While additional development could resolve these
114 issues, it would further aggravate the existing compatibility problem.

115 Another challenge is that some wrappers are infeasible to construct using the above approach. An
116 observation unflattening wrapper, often needed to store observations in flat tensors while retaining the
117 structured format for the model, is one such example. If the flattening wrapper is not the outermost
118 one, the observation space structure required to unflatten the observation is lost. Conversely, if the
119 flattening wrapper is always the final layer, all other wrappers must handle structured observation
120 spaces, thereby adding unnecessary complexity and error-prone code.

121 **4 PufferLib’s Approach**

122 PufferLib aims to handle all the complex data transformations discussed above, returning data in a
123 format compatible with even the most basic reinforcement learning libraries. The system comprises
124 three primary layers: emulation, vectorization, and framework integrations. The ultimate outcome
125 allows users to wrap some of the most intricate reinforcement learning environments available with a
126 single line of code and use a single PyTorch network to train with multiple reinforcement learning
127 frameworks.

128 **4.1 Emulation**

129 This layer forms the core of PufferLib. By applying the aforementioned data transformations, it
130 generates a simple, standard data format, thereby **emulating** the style of simpler environments. Our
131 approach diverges from Gym, PettingZoo, and Supersuit in three significant ways:

- 132 1. PufferLib consists of a single wrapper layer with transformations applied in a fixed sequence.
133 Observations are grouped, then featurized, subsequently flattened, and finally padded and
134 sorted.
- 135 2. It provides fast Cythonized utilities for both flattening and unflattening observations and
136 actions without the issues described earlier.
- 137 3. The wrapper includes substantial precomputation and caching of costly operations.

138 The wrapper class is designed to address all the common difficulties associated with working with
139 complex, multi-agent environments as simply as possible. For context, it totals only around 700 lines
140 of code, which further shrinks excluding the various API usage, input checking exceptions, optional
141 correctness checks, and utility functions. By comparison, the core of PufferLib is shorter than the
142 domain-specific code previously used to support Neural MMO alone. In an ideal world, users would
143 never face uncaught errors in internal libraries. However, as no reinforcement learning library to
144 date has achieved this standard, PufferLib provides a pragmatic solution by offering a simple, single
145 source of failure, as opposed to the potential confusion caused by dozens of conflicting wrappers.

146 4.2 Vectorization

147 Existing vectorization tools built into Gym and PettingZoo lack stable support for multi-agent
148 environments. PufferLib bridges this gap by including a suite of three vectorization tools. These
149 tools leverage the sanitized output format provided by the emulation layer, allowing them to be both
150 performant and simple. Each environment will consistently present the same number of agents, in the
151 same order, with flattened observations. The three vectorization backends are as follows:

- 152 1. **Multiprocessing:** This tool simulates n environments on each of m processes, totaling nm
153 environments, using Python’s native multiprocessing. An additional version, which transfers
154 observations via shared memory, is included. This variant can prove useful for environments
155 demanding high data bandwidth.
- 156 2. **Ray:** This tool, like the multiprocessing one, simulates n environments on each of m
157 processes, using Anyscale’s Ray distributed backend. Although this implementation might
158 be slower for fast environments, it works natively on multi-machine configurations. It also
159 includes a version that transfers observations to the shared Ray memory store instead of
160 directly to processes, which can be faster for specific environment configurations.
- 161 3. **Serial:** This tool simulates all of the environments on a single thread. This setup proves
162 useful for debugging, as it is compatible with breakpoints while maintaining the same
163 API as the previous implementations. Additionally, it is faster for extremely low-latency
164 environments where the overhead of multiprocessing outweighs its benefits.

165 All these backends offer both synchronous and asynchronous APIs, facilitating their use in a buffered
166 setup. In this configuration, the model processes observations for one set of environments while
167 another set of environments processes the previous set of actions. Additionally, all these backends
168 provide hooks for users to shuttle any arbitrary picklable data to the environments. This feature is
169 essential for advanced training methods that need to communicate - for instance, new tasks or maps -
170 with specific environments on remote processes.

171 4.3 Integrations

172 The current release of PufferLib includes support for CleanRL and RLLib, with an extension to
173 Stable Baselines [Raffin et al., 2021] projected for the forthcoming minor versions. It also includes a
174 customized version of CleanRL’s PPO implementation with additional performance optimizations
175 enabled by PufferLib, such as asynchronous but on-policy sampling and improved handling of
176 environments with variable numbers of agents. Owing to the consistent and standard format defined
177 by the emulation layer, even for complex environments, it is relatively straightforward to employ the
178 same PyTorch network across different framework APIs. PufferLib introduces an entirely optional
179 PyTorch base class that separates the *forward()* function into two parts: *encode_observations* and
180 *decode_actions*. Functions preceding a recurrent cell are categorized under the encoding function,
181 and those succeeding it are under the decoding function. This division is implemented because
182 the handling of recurrence is often the most challenging difference among various frameworks. In
183 addition, the mishandling of data reshaping in the recurrent cell is a common source of implementation
184 bugs. We provide additional checks to mitigate this risk. On top of this API, PufferLib constructs
185 a small, per-framework wrapper, which activates the user-specified recurrent cell according to the
186 specific requirements of the given framework. This approach may be expanded to include transformers
187 in the future, although most RL frameworks currently lack support for this.

188 5 Materials Available for Release

189 The public version of PufferLib (version 0.4) is accessible at [pufferaio.github.io](https://github.com/pufferaio/pufferlib). Version 0.5 is planned
190 for release by the end of the year and will include additional framework support. User testing greatly
191 accelerates progress, and the exposure from publication would significantly benefit this work. We
192 currently have the following materials ready for release:

- 193 • Simple documentation and demos for CleanRL and RLlib with Neural MMO available on
194 the website mentioned above.
- 195 • Built-in support and testing for Atari Bellemare et al. [2012], Butterfly (part of PettingZoo),
196 Classic Control (part of Gym), Crafter Hafner [2021], MAgent Zheng et al. [2017], Mi-
197 croRTS [Huang et al., 2021b], Nethack [Küttler et al., 2020], Neural MMO [Suarez et al.,
198 2021], Griddly [Bamford et al., 2020], and partial support for and SMAC [Samvelyan et al.,
199 2019], DM Lab [Beattie et al., 2016], DM Control [Tassa et al., 2018], ProcGen [Cobbe
200 et al., 2019], and MineRL [Guss et al., 2019]. Most of these are one-line wrappers that
201 primarily depend on ensuring compatibility of dependency versions. These are also included
202 in our correctness tests.
- 203 • A Docker container, fondly referred to as PufferTank, that comes pre-built with PufferLib
204 and all of the above environments pre-installed. We have done additional versioning work to
205 maximize the number of common environment that may be installed simultaneously without
206 conflicts.
- 207 • Baselines on the 6 original Atari environments from DQN [Mnih et al., 2013], sanity-checked
208 against CleanRL’s vanilla implementation.
- 209 • A community Discord server with 120 members, offering easy access to support.

210 This version further includes an advanced set of correctness tests that reconstruct the original
211 environment data format from the final version postprocessed by PufferLib. This has aided us in
212 identifying several dozen minor bugs in our development builds. PufferLib is also being utilized
213 in the upcoming Neural MMO competition, enabling much simpler baseline code than would be
214 achievable without it.

215 6 Limitations

216 The most significant limitations of the current release of PufferLib include

- 217 1. No support for Gymnasium. This is already corrected in the 0.5 dev branch.
- 218 2. No support for heterogenous observation and action spaces. These are difficult to process
219 efficiently in a vectorized manner. We have a potential workaround scheduled for a future
220 version.
- 221 3. No support for continuous action spaces. This will be supported with a medium amount of
222 development effort in future versions.
- 223 4. Environments must define a maximum number of agents that fits in memory. Additionally,
224 agents may not respawn. The former is a fundamental limitation of the PettingZoo API. The
225 latter may be supported in with a small amount of development effort.

226 Additionally, as the first publication release of a new framework, we are heavily reliant upon growing
227 a user base to ensure the stability of our tools. We run a battery of correctness tests and verify training
228 performance on Atari in each new release, but subtle bugs have occasionally slipped through during
229 development.

230 7 Conclusion

231 This paper introduces PufferLib, a versatile tool that simplifies working with both single and multi-
232 agent reinforcement learning environments. By providing a consistent data format and handling
233 complex transformations, PufferLib allows researchers to focus on scaling their work to more
234 cognitively interesting environments rather than dealing with finicky compatibility details. Its built-in
235 support for a wide variety of environments, coupled with its scalability and compatibility with popular
236 RL frameworks, makes PufferLib a comprehensive solution for reinforcement learning tasks. We
237 welcome the open-source community to use and contribute to PufferLib, and we anticipate that its
238 ongoing development and integration will lower barriers to reinforcement learning research.

239 References

- 240 Chris Bamford, Shengyi Huang, and Simon M. Lucas. Griddly: A platform for AI research in games.
241 *CoRR*, abs/2011.06363, 2020. URL <https://arxiv.org/abs/2011.06363>.
- 242 Charles Beattie, Joel Z. Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler,
243 Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, Julian Schrittwieser, Keith Anderson,
244 Sarah York, Max Cant, Adam Cain, Adrian Bolton, Stephen Gaffney, Helen King, Demis Hassabis,
245 Shane Legg, and Stig Petersen. Deepmind lab, 2016.
- 246 Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning en-
247 vironment: An evaluation platform for general agents. *CoRR*, abs/1207.4708, 2012. URL
248 <http://arxiv.org/abs/1207.4708>.
- 249 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and
250 Wojciech Zaremba. Openai gym. *CoRR*, abs/1606.01540, 2016. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1606.01540)
251 [1606.01540](http://arxiv.org/abs/1606.01540).
- 252 Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation
253 to benchmark reinforcement learning. *CoRR*, abs/1912.01588, 2019. URL [http://arxiv.org/](http://arxiv.org/abs/1912.01588)
254 [abs/1912.01588](http://arxiv.org/abs/1912.01588).
- 255 William H. Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codell, Manuela Veloso,
256 and Ruslan Salakhutdinov. Minerl: A large-scale dataset of minecraft demonstrations. *CoRR*,
257 abs/1907.13440, 2019. URL <http://arxiv.org/abs/1907.13440>.
- 258 Danijar Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*,
259 2021.
- 260 Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, and Jeff Braga. Cleanrl: High-quality
261 single-file implementations of deep reinforcement learning algorithms. *CoRR*, abs/2111.08819,
262 2021a. URL <https://arxiv.org/abs/2111.08819>.
- 263 Shengyi Huang, Santiago Ontañón, Chris Bamford, and Lukasz Grela. Gym- μ rts: Toward affordable
264 full game real-time strategy games research with deep reinforcement learning. In *2021 IEEE*
265 *Conference on Games (CoG), Copenhagen, Denmark, August 17-20, 2021*, pages 1–8. IEEE,
266 2021b. doi: 10.1109/CoG52621.2021.9619076. URL [https://doi.org/10.1109/CoG52621.](https://doi.org/10.1109/CoG52621.2021.9619076)
267 [2021.9619076](https://doi.org/10.1109/CoG52621.2021.9619076).
- 268 Heinrich Küttler, Nantas Nardelli, Alexander H. Miller, Roberta Raileanu, Marco Selvatici, Edward
269 Grefenstette, and Tim Rocktäschel. The nethack learning environment, 2020.
- 270 Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Joseph Gonzalez, Ken
271 Goldberg, and Ion Stoica. Ray rllib: A composable and scalable reinforcement learning library.
272 *CoRR*, abs/1712.09381, 2017. URL <http://arxiv.org/abs/1712.09381>.
- 273 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
274 Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*,
275 abs/1312.5602, 2013. URL <http://arxiv.org/abs/1312.5602>.
- 276 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
277 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine*
278 *Learning Research*, 22(268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.
- 279 Mikayel Samvelyan, Tabish Rashid, Christian Schröder de Witt, Gregory Farquhar, Nantas Nardelli,
280 Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob N. Foerster, and Shimon Whiteson.
281 The starcraft multi-agent challenge. *CoRR*, abs/1902.04043, 2019. URL [http://arxiv.org/](http://arxiv.org/abs/1902.04043)
282 [abs/1902.04043](http://arxiv.org/abs/1902.04043).

- 283 Joseph Suarez, Yilun Du, Clare Zhu, Igor Mordatch, and Phillip Isola. The neural mmo plat-
 284 form for massively multiagent research. In J. Vanschoren and S. Yeung, editors, *Proceed-*
 285 *ings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, vol-
 286 ume 1, 2021. URL [https://datasets-benchmarks-proceedings.neurips.cc/paper/](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/44f683a84163b3523afe57c2e008bc8c-Paper-round1.pdf)
 287 [2021/file/44f683a84163b3523afe57c2e008bc8c-Paper-round1.pdf](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/44f683a84163b3523afe57c2e008bc8c-Paper-round1.pdf).
- 288 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden,
 289 Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Riedmiller.
 290 Deepmind control suite, 2018.
- 291 Justin K. Terry, Benjamin Black, and Ananth Hari. Supersuit: Simple microwrappers for reinforce-
 292 ment learning environments. *CoRR*, abs/2008.08932, 2020a. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2008.08932)
 293 [2008.08932](https://arxiv.org/abs/2008.08932).
- 294 Justin K. Terry, Benjamin Black, Ananth Hari, Luis S. Santos, Clemens Dieffendahl, Niall L. Williams,
 295 Yashas Lokesh, Caroline Horsch, and Praveen Ravi. Pettingzoo: Gym for multi-agent reinforcement
 296 learning. *CoRR*, abs/2009.14471, 2020b. URL <https://arxiv.org/abs/2009.14471>.
- 297 Lianmin Zheng, Jiacheng Yang, Han Cai, Weinan Zhang, Jun Wang, and Yong Yu. Magent: A
 298 many-agent reinforcement learning platform for artificial collective intelligence, 2017.

299 Checklist

- 300 1. For all authors...
- 301 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 302 contributions and scope? [Yes] We claim only a release of the platform and it’s basic
 303 capabilities, which may be verified from downloading the library.
- 304 (b) Did you describe the limitations of your work? [Yes] See Limitations
- 305 (c) Did you discuss any potential negative societal impacts of your work? [No] This is a
 306 release of tools for academic research
- 307 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 308 them? [Yes]
- 309 2. If you are including theoretical results...
- 310 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 311 (b) Did you include complete proofs of all theoretical results? [N/A]
- 312 3. If you ran experiments (e.g. for benchmarks)...
- 313 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 314 mental results (either in the supplemental material or as a URL)? [Yes] Included in the
 315 base repository, not the package itself
- 316 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 317 were chosen)? [Yes] We used the default hyperparameters of the frameworks
- 318 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 319 ments multiple times)? [No] These experiments were run only as correctness tests to
 320 verify similarity to base CleanRL etc.
- 321 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 322 of GPUs, internal cluster, or cloud provider)? [No] We used a single T40 and 4 cores
 323 for Atari baselines, run for a few days. Given that our work is tooling, this did not seem
 324 relevant to include in the main text.
- 325 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 326 (a) If your work uses existing assets, did you cite the creators? [Yes] Attribution for the
 327 logo and design is provided on the main page

- 328 (b) Did you mention the license of the assets? [Yes] The release (i.e. everything but the
329 logo) is MIT licensed. Copyright for the logo is owned by the author.
- 330 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
331 pufferaai.github.io
- 332 (d) Did you discuss whether and how consent was obtained from people whose data you're
333 using/curating? [N/A] No such data
- 334 (e) Did you discuss whether the data you are using/curating contains personally identifiable
335 information or offensive content? [N/A] No such data
- 336 5. If you used crowdsourcing or conducted research with human subjects...
- 337 (a) Did you include the full text of instructions given to participants and screenshots, if
338 applicable? [N/A] No crowdsourcing
- 339 (b) Did you describe any potential participant risks, with links to Institutional Review
340 Board (IRB) approvals, if applicable? [N/A]
- 341 (c) Did you include the estimated hourly wage paid to participants and the total amount
342 spent on participant compensation? [N/A]