
PufferLib: Making Reinforcement Learning Libraries and Environments Play Nice

Joseph Suarez

Abstract

1 Reinforcement learning (RL) frameworks often falter in complex environments
2 due to inherent simplifying assumptions. This gap necessitates labor-intensive and
3 error-prone intermediate conversion layers, limiting the applicability of RL as a
4 whole. To address this challenge, we introduce PufferLib, a novel middleware
5 solution. PufferLib transforms complex environments into a broadly compatible,
6 vectorized format, eliminating the need for bespoke conversion layers and enabling
7 more rigorous testing. Users interact with PufferLib through concise bindings, sig-
8 nificantly reducing the technical overhead. We release PufferLib’s complete source
9 code under the MIT license, a pip module, a containerized setup, comprehensive
10 documentation, and example integrations. We also maintain a community Discord
11 channel to facilitate support and discussion.

12 1 Background and Introduction

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

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

26 In response, we introduce PufferLib, a novel middleware solution bridging complex environments
27 and reinforcement learning libraries, effectively mitigating the integration challenges. PufferLib
28 decouples one layer of RL’s unique complexities, making the remaining challenges more manageable
29 and fostering more rapid progress in the field. Other existing solutions such as Gym [Brockman et al.,
30 2016], PettingZoo [Terry et al., 2020b], and SuperSuit [Terry et al., 2020a] aim to define standard
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 for use with popular
35 reinforcement learning libraries, such as CleanRL [Huang et al., 2021a] and RLlib [Liang et al.,

The PufferLib System Architecture for Broad Environment Compatibility

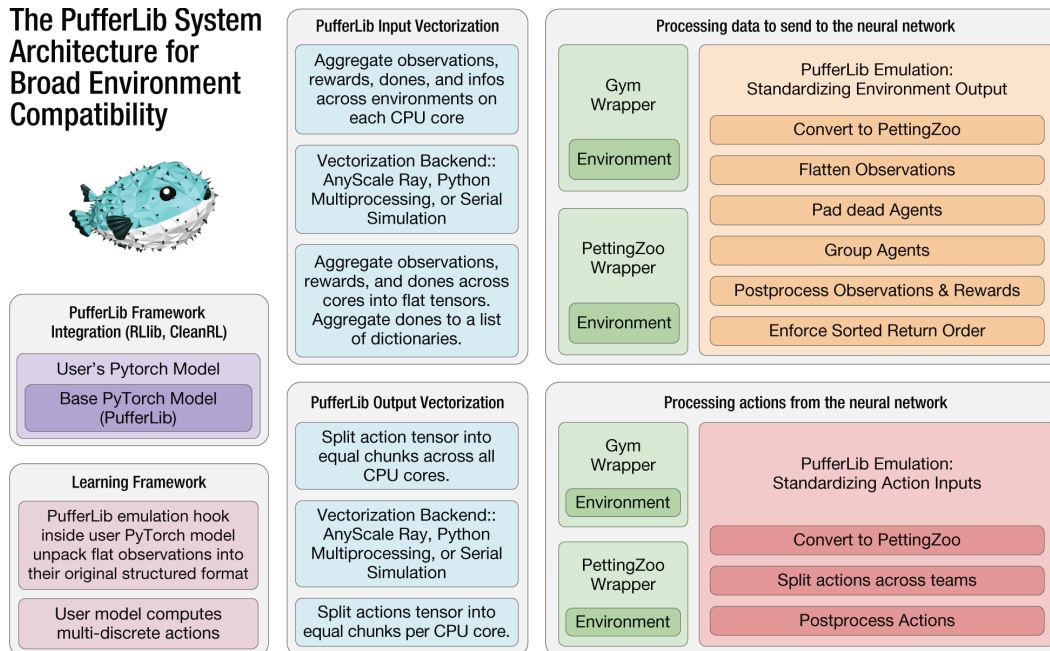


Figure 1: Detailed but non-comprehensive illustration of the PufferLib system architecture, comprising emulation, vectorization, and learning framework integrations. The orange emulation block demonstrate how PufferLib receives and processes environment data. The red emulation block demonstrates how 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.

36 2017]. It natively supports multi-agent and variable-agent environments and addresses common
 37 complexities that include batching structured observations and actions, featurizing observations,
 38 shaping rewards, and grouping agents into teams. PufferLib is also designed for extensibility and is
 39 capable of supporting new learning libraries with a complete feature set in typically about a hundred
 40 lines of code.

41 2 Problem Statement

42 To thoroughly ground our work, we will walk through the intricacies of the transformations that
 43 reinforcement learning data must undergo, and demonstrate the shortcomings of existing approaches.
 44 Specifically, we will trace the required transformations from simulation onset to data processing by
 45 the initial model layer, and from action computation to the point when those actions influence the
 46 environment.

47 We will use Neural MMO [Suarez et al., 2021], a Gym and PettingZoo-compliant environment, as
 48 our guiding example. This environment, encapsulating many complexities common to advanced
 49 environments, features 128 agents competing to complete tasks in a procedurally-generated world.
 50 It provides agents with rich, structured observations of their surroundings and a hierarchical action
 51 space for interactions.

52 The environment initialization starts with a configuration file and a reset to yield an initial set of
 53 observations. This results in a dictionary of 128 individual observations, each of which is a structured
 54 dictionary housing differently-shaped tensors related to various aspects of the observation. As a part
 55 of the environment’s standard training setting, these agents are grouped into teams of 8. Each team
 56 observation is then processed by a featurizer to yield a single structured observation, aggregating

57 information from across the team’s agents. Subsequently, this observation must be batched for model
58 usage.

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

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

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

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

90 **3 Related Tools**

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

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

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

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

119 **4 PufferLib’s Approach**

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

126 **4.1 Emulation**

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

- 130 1. PufferLib consists of a single wrapper layer with transformations applied in a fixed sequence.
131 Observations are grouped, then featurized, subsequently flattened, and finally padded and
132 sorted.
- 133 2. It provides utilities for both flattening and unflattening observations and actions without the
134 issues described earlier.
- 135 3. The wrapper class is procedurally generated using data scoped from a dummy instance of
136 the unwrapped environment, enabling the static precomputation of a few costly operations.

137 The emulation layer starts with a Binding object. Users can instantiate a binding from a Gym or
138 PettingZoo environment class, instance, or creation function. They can supply several arguments to
139 the Binding object, including a custom postprocessor for features, actions, and rewards, choices about
140 flattening observation and action spaces, whether to pad to a constant number of agents, whether
141 to truncate environment simulation at a set number of steps, etc. The Binding class creates or uses
142 the provided environment instance and resets it to yield an initial observation. This observation,
143 alongside the provided binding arguments, is used to create a wrapper class for the environment. The
144 significance of this process is that it allows the initial observation to be statically scoped into the
145 wrapper class. The Binding then offers access to the wrapper class with no intermediate layer.

146 The wrapper class is designed to address all the common difficulties associated with working with
147 complex, multi-agent environments as simply as possible. For context, it totals only around 800 lines
148 of code, which further shrinks excluding the various API usage, input checking exceptions, optional
149 correctness checks, and utility functions. By comparison, the core of PufferLib is shorter than the
150 domain-specific code previously used to support Neural MMO alone. In an ideal world, users would
151 never face uncaught errors in internal libraries. However, as no reinforcement learning library to

152 date has achieved this standard, PufferLib provides a pragmatic solution by offering a simple, single
153 source of failure, as opposed to the potential confusion caused by dozens of conflicting wrappers.

154 4.2 Vectorization

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

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

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

179 4.3 Integrations

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

193 5 Materials Available for Release

194 The public version of PufferLib (version 0.3) is accessible at [pufferai.github.io](https://github.com/pufferai/pufferlib). Version 0.4 is planned
195 for release by the end of the summer and will include additional framework support. User testing
196 greatly accelerates progress, and the exposure from publication would significantly benefit this work.
197 We currently have the following materials ready for release:

- 198 • Simple documentation and demos for CleanRL and RLib with Neural MMO available on
199 the website mentioned above.
- 200 • Built-in support and testing for Atari Bellemare et al. [2012], Butterfly (part of PettingZoo),
201 Classic Control (part of Gym), Crafter Hafner [2021], MAgent Zheng et al. [2017], Mi-
202 croRTS [Huang et al., 2021b], Nethack [Küttler et al., 2020], Neural MMO [Suarez et al.,
203 2021], and SMAC [Samvelyan et al., 2019] with partial support for Griddly [Bamford et al.,
204 2020] and planned extensions to DM Lab [Beattie et al., 2016], DM Control [Tassa et al.,
205 2018], ProcGen [Cobbe et al., 2019], and MineRL [Guss et al., 2019]. Most of these are
206 one-line bindings that primarily depend on ensuring compatibility of dependency versions.
207 These are also included in our correctness tests.
- 208 • A Docker container, fondly referred to as PufferTank, that comes pre-built with PufferLib
209 and all of the above environments pre-installed.
- 210 • Baselines on the 6 original Atari environments from DQN [Mnih et al., 2013], sanity-checked
211 against CleanRL’s vanilla implementation.
- 212 • A community Discord server with 100 members, offering easy access to support.

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

218 6 Limitations

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

- 220 1. No support for heterogenous observation and action spaces. These are difficult to process
221 efficiently in a vectorized manner.
- 222 2. No support for continuous action spaces. This may be supported with a medium amount of
223 development effort in future versions.
- 224 3. Environments must define a maximum number of agents that fits in memory. Additionally,
225 agents may not respawn. The former is a fundamental limitation of the underlying PettingZoo
226 binding. The latter may be supported in a future version with a small amount of development
227 effort.

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

232 7 Conclusion

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

241 References

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

285 Joseph Suarez, Yilun Du, Clare Zhu, Igor Mordatch, and Phillip Isola. The neural mmo plat-
 286 form for massively multiagent research. In J. Vanschoren and S. Yeung, editors, *Proceed-*
 287 *ings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, vol-
 288 ume 1, 2021. URL [https://datasets-benchmarks-proceedings.neurips.cc/paper/](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/44f683a84163b3523afe57c2e008bc8c-Paper-round1.pdf)
 289 [2021/file/44f683a84163b3523afe57c2e008bc8c-Paper-round1.pdf](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/44f683a84163b3523afe57c2e008bc8c-Paper-round1.pdf).

290 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden,
 291 Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Riedmiller.
 292 Deepmind control suite, 2018.

293 Justin K. Terry, Benjamin Black, and Ananth Hari. Supersuit: Simple microwrappers for reinforce-
 294 ment learning environments. *CoRR*, abs/2008.08932, 2020a. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2008.08932)
 295 [2008.08932](https://arxiv.org/abs/2008.08932).

296 Justin K. Terry, Benjamin Black, Ananth Hari, Luis S. Santos, Clemens Dieffendahl, Niall L. Williams,
 297 Yashas Lokesh, Caroline Horsch, and Praveen Ravi. Pettingzoo: Gym for multi-agent reinforcement
 298 learning. *CoRR*, abs/2009.14471, 2020b. URL <https://arxiv.org/abs/2009.14471>.

299 Lianmin Zheng, Jiacheng Yang, Han Cai, Weinan Zhang, Jun Wang, and Yong Yu. Magent: A
 300 many-agent reinforcement learning platform for artificial collective intelligence, 2017.

301 Checklist

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

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