
Syllabus: Curriculum Learning Made Easy

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Curriculum learning has been a quiet yet crucial component of many of the high-
2 profile successes of reinforcement learning. Despite this, none of the major re-
3 inforcement learning libraries support curriculum learning or include curriculum
4 learning algorithms. Curriculum learning methods can provide general and comple-
5 mentary improvements to RL algorithms, but they often require significant,
6 complex changes to agent training code. We introduce Syllabus, a library for
7 training RL agents with curriculum learning, as a solution to this problem. Syllabus
8 provides a universal API for implementing curriculum learning algorithms, a collec-
9 tion of implementations of popular curriculum learning methods, and infrastructure
10 for easily integrating them into existing distributed RL code. Syllabus provides a
11 clean API for each of the complex components of these methods, dramatically sim-
12 plifying the process for designing new algorithms or applying existing algorithms
13 to new environments. Syllabus also manages the multiprocessing communication
14 required for curriculum learning, alleviating one of the key practical challenges of
15 using these algorithms. We hope Syllabus will improve the process of developing
16 and applying curriculum learning algorithms, and encourage widespread adaptation
17 of curriculum learning.

18 1 Introduction

19 Curricula have been a core component of many of the successes of reinforcement learning. AlphaGo
20 [Silver et al., 2016] was trained with self-play, AlphaStar used a novel league training method to
21 achieve grandmaster level play in Starcraft II [Vinyals et al., 2019], and GT Sophy [Wurman et al.,
22 2022] was recently trained to outrace professionals in Gran Turismo with a manually designed
23 curriculum. Curriculum learning methods, both automatic and manual, fit naturally into the RL
24 framework by modifying the distribution of states that an agent experiences. Both in theory and
25 in practice, these methods often provide orthogonal benefits to reinforcement learning algorithms
26 and can be implemented on top of most RL algorithms with minimal restrictions. Despite the
27 near ubiquity of curriculum learning in successful applications of reinforcement learning, there
28 is relatively little support for these methods in standard RL tools, largely due the complexity of
29 curriculum learning code. Curriculum learning methods may modify many parts of the training
30 process, such as the environment initialization and reward functions, and might even introduce new
31 neural networks to train. Automatic curriculum learning methods may update their internal state
32 based on information from the training process or from the environment which typically run in
33 separate processes. In practice, this entanglement of curriculum learning and RL code makes it
34 challenging to take curriculum learning source code and apply it to a new project with pre-existing
35 RL code.

36 Syllabus addresses this challenge in a number of ways. Even the task of structuring the software
37 architecture for curriculum learning can be daunting. Syllabus clearly designates responsibilities to
38 a separate Curriculum class, and implements changes to the environment through a Task Wrapper,

39 limiting curriculum code to two locations in the code base. Within the Task Wrapper, users can also
40 define a Task Space for their environment that represents the full space of tasks for an environment.
41 These APIs are described in more detail in section 4. We validate our API by implementing several
42 popular curriculum learning methods, demonstrating that it is sufficiently general. Some methods
43 like league training [Vinyals et al., 2019] may be outside the range of Syllabus’s Curriculum API, but
44 Syllabus makes it easy to combine one of its methods with your own custom curriculum code.

45 **2 Background**

46 Curriculum Learning has been studied in the context of deep learning for many years [Bengio et al.,
47 2009, Elman, 1993]. More recently, it has been used to improve the performance of reinforcement
48 learning agents. Curriculum learning encompasses a wide range of methods targeted at changing the
49 distribution of data used to train an agent. The goal is to increase the final performance or sample
50 efficiency of RL agents on a single environment or range of tasks by sampling tasks that provide the
51 maximum learning value. Narvekar et al. [2020] presents a more thorough taxonomy and survey of
52 existing curriculum learning methods.

53 Many diverse methods fall under the broad definition of curriculum learning. Exploration bonuses
54 like curiosity [Pathak et al., 2017] or [Bellemare et al., 2016, Taiga et al., 2021, Henaff et al., 2022]
55 induce a curriculum by incentivizing the agent to explore unseen section of the state space. Methods
56 like self-play or league-play [Vinyals et al., 2019] create an implicit curriculum by training the
57 opponents in a multiplayer game [Leibo et al., 2019]. Progressively more capable opponents lead
58 to progressively more difficult tasks for an agent. Most task-based methods follow the general rule
59 of proposing tasks that are hard yet solvable or tasks which the agent is recently performing well
60 on [Graves et al., 2017, Kanitscheider et al., 2021]. In general, curriculum methods try to create a
61 sequence of increasingly complex or difficult tasks.

62 Unsupervised Environment Design (UED) is another paradigm for curriculum learning proposed
63 by Dennis et al. [2020]. They differentiate UED as a framework in which environments have
64 unspecified parameters, thereby forming an underspecified MDP. Those parameters are used to
65 produce a distribution of solvable tasks, and the goal of a UED method is to produce a policy that
66 generalizes across a large set of tasks.

67 **3 Design Philosophy**

68 Syllabus was designed to solve a number of challenges in automatic curriculum learning. Our main
69 goal is to simplify the process of testing and developing new curriculum learning methods. Syllabus
70 is also built to manage the multiprocessing communication required for curriculum learning, which
71 can be challenging to combine with existing distributed reinforcement learning.

72 These objectives motivated the following key points of our design philosophy:

- 73 1. Syllabus should be agnostic to the choice of reinforcement learning framework.
- 74 2. Syllabus should be as general as possible to support many popular curriculum learning
75 methods.
- 76 3. Using Syllabus should require minimal changes to existing RL code.
- 77 4. If a method requires more complex code, complexity of the code should scale with the
78 complexity of the curriculum learning method. Simple methods should remain simple to
79 use.
- 80 5. Single-file implementations of individual algorithms.

81 The first point motivates many of the implementation choices in Syllabus which may seem odd
82 in isolation. It requires that we honor the Gym [Brockman et al., 2016] environment API and
83 write systems that work well with the many different forms of multiprocessing used throughout
84 reinforcement learning libraries.

85 The second, third, and fourth point pertain to the Curriculum API that Syllabus defines. These
86 goals often conflict due to the wide range of methods that need to be supported under a single API.
87 Some curricula may require metrics from the training process while others only utilize rewards

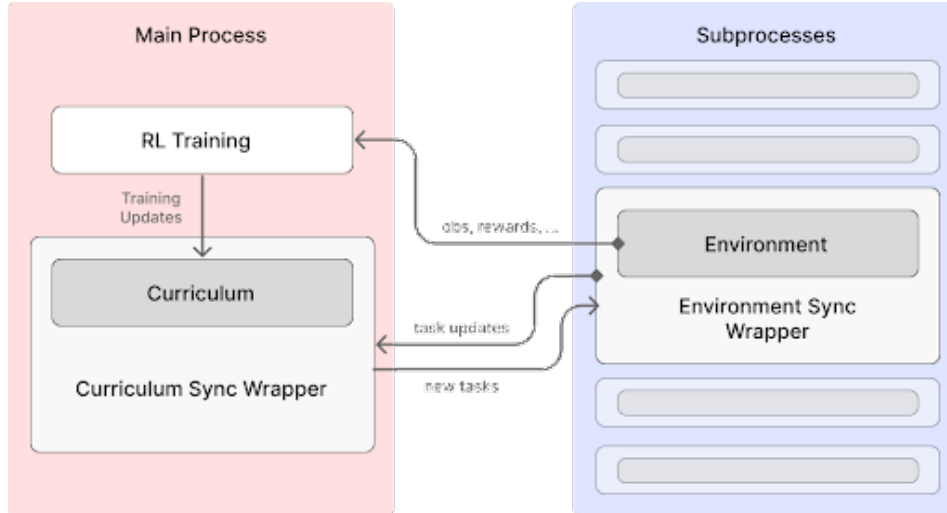


Figure 1: Syllabus with a standard distributed RL training setup.

88 from the environments. These components typically run in separate processes and therefore require
 89 very different interfaces. In line with the fourth design goal, when one method requires additional
 90 changes to the training code, we prefer to have a heterogenous API for different methods rather than
 91 complicate the API for all methods. To reduce confusion, we thoroughly document these changes
 92 and raise warnings for common user errors resulting from those differences. We explain our approach
 93 to multiprocessing and balancing these design challenges more in subsection 5.2.

94 Finally, we choose to use single-file or few-file implementations of curriculum methods inspired by
 95 the success of CleanRL [Huang et al., 2022], a popular RL library which strictly uses single-file
 96 implementations to simplify research engineering.

97 4 Syllabus APIs

98 Syllabus defines APIs for the key components of curriculum learning; the Curriculum API, the Task
 99 Space API, and the Task Wrapper API. These are the main components that users will need to modify
 100 or interact with to develop or apply curriculum learning methods, so they are each designed to be
 101 both simple and general enough to support future use cases.

102 4.1 Curriculum API

103 In Syllabus, a Curriculum is responsible for maintaining a distribution over the task space and
 104 implementing a sampling function for selecting tasks. Automatic curriculum learning methods
 105 require feedback from the RL process to update their sampling distribution. Syllabus provides
 106 multiple optional methods which a Curriculum can implement to receive updates from different
 107 sources. It can either be updated directly by the main training process, or it can automatically
 108 receive updates from the environments through Syllabus’s multiprocessing infrastructure. The main
 109 components of the API are shown in Figure 2.

110 4.2 Task Space API

111 Most curriculum learning methods operate over a distribution of tasks. The task space API defines
 112 the range and bounds of this distribution, as well as a conversion between a simple task identifier and
 113 the full task definition expected by the environment. This simple identifier can be used to index the
 114 task within the task space which allows us to define a subset of validation tasks.

115 One of the design challenges of applying curriculum learning to a new domain is defining the task
 116 space of the environment. In most benchmark environments (like Procgen [Cobbe et al., 2020] or
 117 Minigrid [Chevalier-Boisvert et al., 2023]) task spaces are low-dimensional discrete or continuous

Figure 2: Syllabus’s Curriculum API.

```

1 class Curriculum:
2     """API for defining curricula to interface with Gym environments."""
3
4     @property
5     def num_tasks(self) -> int:
6         """Counts the number of tasks in the task space, if countable."""
7
8     @property
9     def tasks(self) -> List[tuple]:
10        """ List all of the tasks in the task space, if enumerable."""
11
12    def update_task_progress(
13        self, task: Any, progress: Tuple[float, bool]
14    ) -> None:
15        """ Update the curriculum with a task and its progress.
16            Progress is defined by the environment's TaskWrapper. """
17
18    def update_on_step(
19        self, obs: Any, rew: float, done: bool, info: dict
20    ) -> None:
21        """ Update the curriculum with the environment outputs
22            for the most recent step. """
23
24    def update_on_demand(self, metrics: Dict):
25        """ Update the curriculum with arbitrary inputs.
26            Typically used to incorporate gradient or error-based
27            metrics from the training process. """
28
29    def _sample_distribution(self) -> List[float]:
30        """ Returns a sample distribution over the task space.
31            Any curriculum that maintains a true probability distribution
32            should implement this method to retrieve the distribution. """
33
34    def sample(self, k: int = 1) -> List:
35        """ Sample k tasks from the curriculum. """
36

```

118 spaces. In many real environments the space of possible tasks might be a combination of discrete and
119 continuous variables, or a complex predicate system such as in Neural MMO [Suarez et al., 2019].
120 Curriculum learning algorithms typically only support one specific task space representation. The
121 Task Space API is designed to alleviate some of the challenges of defining task spaces and identifying
122 which curriculum methods can be used with a particular task space. It automatically encodes tasks
123 into a Gym Space which defines the full range of tasks.

124 4.3 Task Wrapper API

125 Outside of benchmark environments, the user might need to write code to reinitialize an environment
126 for each task. Syllabus provides a Task Wrapper API to facilitate this configuration process. The main
127 change that this introduces to the standard Gym API is adding a ‘task’ property to the environment
128 and a ‘new_task’ argument to the environment’s ‘reset()’ function. This allows the user to optionally
129 assign a new task to the environment while resetting the environment. If necessary, this API can also
130 be used to add a task system to an environment that does not natively support one. For example, in
131 Figure 3 we use this API to add a simple task system to CartPole which changes the range from which
132 the initial location of the cart is sampled. Finally, the task wrapper can also define a progress metric,
133 a float value in the range [0.0, 1.0] that defines how complete the current task is. In the simplest case,

134 this can be a binary value that is 1.0 when the task is completed and 0.0 while it is not. This allows
135 Syllabus to support learning-progress metrics like the one proposed in Kanitscheider et al. [2021].

136 **5 Implementation**

137 **5.1 Implemented Curriculum Learning Methods**

138 Syllabus provides single-file implementations of popular curriculum learning methods including
139 Domain Randomization, Prioritized Level Replay [Jiang et al., 2021a] and the learning progress
140 curriculum proposed by Kanitscheider et al. [2021], with plans to add many more. Each of these
141 methods are tested with Syllabus’s infrastructure and include warnings for common user errors.

142 **5.2 Multiprocessing Infrastructure**

143 Syllabus’s multiprocessing infrastructure uses a bidirectional sender-receiver model in which the
144 curriculum sends tasks and receives environment outputs, while the environment receives tasks and
145 sends outputs. Environments run the provided task in the next episode and the curriculum can use the
146 outputs to update its task distribution. This communication layer is implemented in two wrappers
147 that add functionality while maintain the same interface. The curriculum synchronization wrapper
148 adds multiprocessing functionality to a Curriculum and an environment synchronization wrapper,
149 similar to the environment wrappers used extensively in reinforcement learning code, adds the same
150 functionality to the environment. You can also call the curriculum directly from the main learner
151 process to update it with training metrics. Crucially, adding Syllabus’s functionality to existing RL
152 training code requires only a few lines of code. Figure 1 shows a diagram of how these components
153 interconnect. These synchronization wrappers automatically send the curriculum environment outputs
154 at each step and episodic updates on task progress. All updates are batched to reduce multiprocessing
155 overhead, and the per-step updates can be disabled to improve performance if they are not used by
156 the chosen curriculum learning method.

157 Most of the user facing Curriculum and environment code follows our design goal of single-file
158 implementations, while the multiprocessing infrastructure is more engineered to ensure stability
159 and reduce the risk of bugs. To guarantee that researchers will not need to spend time reading
160 or debugging this code, Syllabus include thorough integration tests, smoke tests, regression tests,
161 and optimization benchmarks for the entire multiprocessing infrastructure, tested with all of the
162 implemented curriculum learning methods.

163 **5.3 Integration**

164 Following our goal of requiring minimal code changes, the example in Figure 3 shows how easy it is
165 to add a simulated annealing curriculum to RLLib training code. In addition, because each method
166 is implemented with the same API, replacing this curriculum with another only requires a single
167 change. Simply initialize a different curriculum (line 28 in Figure 3), and it will work as expected for
168 most methods. For the few methods that require additional changes, warnings will be raised if those
169 changes are not made. For example, if Prioritized Level Replay does not receive any TD errors from
170 the training process after a full batch of environment experience, it will raise an error and direct
171 users to documentation for adding the missing code.

172 **5.4 Optimization**

173 As a consequence of the choice to use a separate multiprocessing system from the RL training
174 loop, Syllabus incurs some unavoidable computational costs. Specifically, receiving and sending
175 information in the environments decreases the effective steps per second of each environment, while
176 sampling and sending tasks in the actor process increases the computational load on the main process.
177 We perform experiments on the NetHack Learning Environment [Küttler et al., 2020] to demonstrate
178 the effect of this choice on overall steps per second. We evaluate with Domain Randomization, a
179 computationally lightweight method, to isolate the impact of our multiprocessing infrastructure. The
180 results are show in Table 1 for experiments run with 128 environments and 50 episodes on a 32 core
181 Intel i9-13950HX. We test Syllabus using both Python’s native multiprocessing package and Ray as
182 the backend. Syllabus allows you to disable per-step updates for ACL methods that do not require

Figure 3: Adding curriculum learning with Syllabus to RLLib training code with just a few lines of code.

```

1 import gym
2 from ray.tune.registry import register_env
3 from ray import tune
4 from gym.spaces import Box
5 from .task_wrappers import CartPoleTaskWrapper
6
7 +from syllabus.core import RaySyncWrapper, make_ray_curriculum
8 +from syllabus.curricula import SimpleBoxCurriculum
9 +from syllabus.task_space import TaskSpace
10
11 if __name__ == "__main__":
12 + # Define a task space
13 + task_space = TaskSpace(Box(-0.3, 0.3, shape=(2,)), [])
14
15 def env_creator(config):
16     env = gym.make("CartPole-v1")
17     # Wrap the environment to change tasks on reset()
18     env = CartPoleTaskWrapper(env)
19 + # Add environment sync wrapper
20 + env = RaySyncWrapper(
21 +     env, default_task=(-0.02, 0.02), task_space=task_space
22 + )
23     return env
24
25 register_env("task_cartpole", env_creator)
26
27 + # Create the curriculum
28 + curriculum = SimpleBoxCurriculum(task_space)
29 + # Add the curriculum sync wrapper
30 + curriculum = make_ray_curriculum(curriculum)
31
32 config = {
33     "env": "task_cartpole",
34     "num_gpus": 1,
35     "num_workers": 8,
36     "framework": "torch",
37 }
38
39 tuner = tune.Tuner("APEX", param_space=config)
40 results = tuner.fit()
41

```

183 them, instead only sending updates and new tasks at the end of each episode. We show results for
184 both of these scenarios. Note that the NLE is an extremely fast environment, and often RL training
185 with larger architectures is bottle-necked by policy optimization rather than environment iteration
186 time. We expect Syllabus’s impact on performance (as a percentage of total computation) to be much
187 lower for more computationally intensive environments and when combined with RL training.

188 5.5 Implemented Curriculum Learning Methods

189 Curriculum learning algorithms vary as much as RL algorithms in complexity. The simplest methods
190 like Domain Randomization never update their sampling distributions and have simple logic for
191 sampling tasks. The most complex methods, such as league training, might need to train entire
192 new agents and maintain many sets of network weights. Syllabus aims to eventually support all

Table 1: Syllabus Performance Costs

Multiprocessing	NLE	Syllabus (Episodic Updates)	Syllabus (Step Updates)
Native Python	58.5s	63.0s (+7.7%)	93.0s (+59.0%)
Ray	71.12s (100%)	84.7s (+19.1%)	134.2s (+88.7%)

193 curriculum learning paradigms, but it is currently limited to a subset of paradigms that minimally
 194 alter the training code. Table x shows some of the key implementation requirements of curriculum
 195 learning methods and how they correspond to some popular methods.

196 6 Related Work

197 The closest work to Syllabus is RLLib [Liang et al., 2018], which is a large library of reinforcement
 198 learning algorithms. RLLib allows you to set the task of an environment, which is similar to Syllabus’s
 199 Task Environment API, and it manages all of the multiprocessing through Ray. RLLib does not
 200 implement any curriculum learning methods or provide utilities for updating a curriculum using
 201 metrics from the environment or training process, though it would be possible to implement these
 202 through a combination of callbacks and Ray Actors.

203 Another related library is the Dual Curriculum Design library [Jiang et al., 2022], which incorporates
 204 multiple Unsupervised Environment Design Methods in a single repository. It includes implementa-
 205 tions of PLR [Jiang et al., 2021a], PAIRED [Dennis et al., 2020, Mediratta et al., 2023], Robust PLR,
 206 REPAIRED [Jiang et al., 2021b], and ACCEL [Parker-Holder et al., 2022]. Unlike Syllabus, which
 207 provides tools for adding curriculum learning to existing RL code, the DCD library encourages users
 208 to iterate on top of it by adding new RL algorithms or environments.

209 7 Conclusion

210 Curriculum learning is a sub-field of reinforcement learning that is under-utilized and under-
 211 researched relative to its impact on well-known RL applications and empirical benefits. We believe
 212 that this is largely due to the practical and engineering challenges associated with these methods.
 213 Syllabus is open-sourced on GitHub with a complete documentation website, and can be installed as
 214 a pip package. We are continuing to actively develop Syllabus and hope that it will help to encourage
 215 the more widespread use of curriculum learning.

216 References

- 217 David Silver, Aja Huang, Christopher Maddison, Arthur Guez, Laurent Sifre, George Driessche,
 218 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman,
 219 Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine
 220 Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with
 221 deep neural networks and tree search. *Nature*, 529:484–489, 01 2016. doi: 10.1038/nature16961.
- 222 Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Juny-
 223 oung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan
 224 Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, L. Sifre, Trevor Cai, John P. Agapiou, Max
 225 Jaderberg, Alexander Sasha Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David
 226 Budden, Yury Sulsky, James Molloy, Tom Le Paine, Caglar Gulcehre, Ziyun Wang, Tobias Pfaff,
 227 Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom
 228 Schaul, Timothy P. Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver.
 229 Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575:350 – 354,
 230 2019. URL <https://api.semanticscholar.org/CorpusID:204972004>.
- 231 Peter Wurman, Samuel Barrett, Kenta Kawamoto, James MacGlashan, Kaushik Subramanian,
 232 Thomas Walsh, Roberto Capobianco, Alisa Devlic, Franziska Eckert, Florian Fuchs, Leilani
 233 Gilpin, Piyush Khandelwal, Varun Kompella, HaoChih Lin, Patrick MacAlpine, Declan Oller,
 234 Takuma Seno, Craig Sherstan, Michael Thomure, and Hiroaki Kitano. Outracing champion

- 235 gran turismo drivers with deep reinforcement learning. *Nature*, 602:223–228, 02 2022. doi:
236 10.1038/s41586-021-04357-7.
- 237 Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In
238 *Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09*, page
239 41–48, New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605585161.
240 doi: 10.1145/1553374.1553380. URL <https://doi.org/10.1145/1553374.1553380>.
- 241 Jeffrey L. Elman. Learning and development in neural networks: the importance of starting small. *Cog-*
242 *nition*, 48(1):71–99, 1993. ISSN 0010-0277. doi: [https://doi.org/10.1016/0010-0277\(93\)90058-4](https://doi.org/10.1016/0010-0277(93)90058-4).
243 URL <https://www.sciencedirect.com/science/article/pii/0010027793900584>.
- 244 Sanmit Narvekar, Bei Peng, Matteo Leonetti, Jivko Sinapov, Matthew E Taylor, and Peter Stone.
245 Curriculum learning for reinforcement learning domains: A framework and survey. *The Journal of*
246 *Machine Learning Research*, 21(1):7382–7431, 2020.
- 247 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by
248 self-supervised prediction. In *International conference on machine learning*, pages 2778–2787.
249 PMLR, 2017.
- 250 Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.
251 Unifying count-based exploration and intrinsic motivation. *Advances in neural information*
252 *processing systems*, 29, 2016.
- 253 Adrien Ali Taiga, William Fedus, Marlos C Machado, Aaron Courville, and Marc G Belle-
254 mare. On bonus-based exploration methods in the arcade learning environment. *arXiv preprint*
255 *arXiv:2109.11052*, 2021.
- 256 Mikael Henaff, Roberta Raileanu, Minqi Jiang, and Tim Rocktäschel. Exploration via elliptical
257 episodic bonuses. *Advances in Neural Information Processing Systems*, 35:37631–37646, 2022.
- 258 Joel Z Leibo, Edward Hughes, Marc Lanctot, and Thore Graepel. Autocurricula and the emergence
259 of innovation from social interaction: A manifesto for multi-agent intelligence research. *arXiv*
260 *preprint arXiv:1903.00742*, 2019.
- 261 Alex Graves, Marc G Bellemare, Jacob Menick, Remi Munos, and Koray Kavukcuoglu. Automated
262 curriculum learning for neural networks. In *international conference on machine learning*, pages
263 1311–1320. Pmlr, 2017.
- 264 Ingmar Kanitscheider, Joost Huizinga, David Farhi, William Hebgen Guss, Brandon Houghton,
265 Raul Sampedro, Peter Zhokhov, Bowen Baker, Adrien Ecoffet, Jie Tang, et al. Multi-task cur-
266 riculum learning in a complex, visual, hard-exploration domain: Minecraft. *arXiv preprint*
267 *arXiv:2106.14876*, 2021.
- 268 Michael Dennis, Natasha Jaques, Eugene Vinitsky, Alexandre Bayen, Stuart Russell, Andrew Critch,
269 and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised environment
270 design. *Advances in neural information processing systems*, 33:13049–13061, 2020.
- 271 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and
272 Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 273 Shengyi Huang, Rousslan Fernand JulienDossa Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty,
274 Kinal Mehta, and João GM Araújo. Cleanrl: High-quality single-file implementations of deep
275 reinforcement learning algorithms. *The Journal of Machine Learning Research*, 23(1):12585–
276 12602, 2022.
- 277 Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation
278 to benchmark reinforcement learning. In *International conference on machine learning*, pages
279 2048–2056. PMLR, 2020.
- 280 Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem
281 Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular &
282 customizable reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831,
283 2023.

- 284 Joseph Suarez, Yilun Du, Phillip Isola, and Igor Mordatch. Neural mmo: A massively multiagent
285 game environment for training and evaluating intelligent agents. *arXiv preprint arXiv:1903.00784*,
286 2019.
- 287 Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay. In *International*
288 *Conference on Machine Learning*, pages 4940–4950. PMLR, 2021a.
- 289 Heinrich Küttler, Nantas Nardelli, Alexander Miller, Roberta Raileanu, Marco Selvatici, Edward
290 Grefenstette, and Tim Rocktäschel. The nethack learning environment. *Advances in Neural*
291 *Information Processing Systems*, 33:7671–7684, 2020.
- 292 Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph
293 Gonzalez, Michael Jordan, and Ion Stoica. Rllib: Abstractions for distributed reinforcement
294 learning. In *International conference on machine learning*, pages 3053–3062. PMLR, 2018.
- 295 Minqi Jiang, Ishita Mediratta, Mikayel Samvelyan, and Jayden Teoh. Dual curriculum design.
296 <https://github.com/facebookresearch/dcd>, 2022.
- 297 Ishita Mediratta, Minqi Jiang, Jack Parker-Holder, Michael Dennis, Eugene Vinitzky, and Tim
298 Rocktäschel. Stabilizing unsupervised environment design with a learned adversary. *arXiv preprint*
299 *arXiv:2308.10797*, 2023.
- 300 Minqi Jiang, Michael Dennis, Jack Parker-Holder, Jakob Foerster, Edward Grefenstette, and Tim
301 Rocktäschel. Replay-guided adversarial environment design. *Advances in Neural Information*
302 *Processing Systems*, 34:1884–1897, 2021b.
- 303 Jack Parker-Holder, Minqi Jiang, Michael Dennis, Mikayel Samvelyan, Jakob Foerster, Edward
304 Grefenstette, and Tim Rocktäschel. Evolving curricula with regret-based environment design. In
305 *International Conference on Machine Learning*, pages 17473–17498. PMLR, 2022.