

---

# simple\_rl: Reproducible Reinforcement Learning in Python

---

David Abel  
Brown University  
Providence, RI 02903  
david.abel@brown.edu

## Abstract

1     Conducting reinforcement-learning experiments can be a complex and timely process.  
2     A full experimental pipeline will typically consist of a simulation of an environment,  
3     an implementation of one or many learning algorithms, a variety of  
4     additional components designed to facilitate the agent-environment interplay, and  
5     any requisite analysis, plotting, and logging thereof. In light of this complexity,  
6     this paper introduces `simple_rl`<sup>1</sup>, a new open source library for carrying out reinforcement  
7     learning experiments in Python 2 and 3 with a focus on simplicity. The goal of `simple_rl`  
8     is to support seamless, reproducible methods for running reinforcement learning experiments.  
9     This paper gives an overview of the core design philosophy of the package, how it differs  
10    from existing libraries, and showcases its central features.  
11



Figure 1: The core functionality of `simple_rl`: Create agents and an MDP, then run and plot their resulting interactions. Running an experiment also creates an experiment log (stored as a JSON file), which can be used to rerun the exact same experiment, thereby facilitating simple reproduction of results. All practitioners need to do, in theory, is share a copy of the experiment file to someone with the library to ensure result reproduction.

<sup>1</sup>[https://github.com/david-abel/simple\\_rl](https://github.com/david-abel/simple_rl)

---

```

1  from simple_rl.agents import QLearningAgent, RandomAgent
2  from simple_rl.tasks import GridWorldMDP
3  from simple_rl.run_experiments import run_agents_on_mdp
4
5  # Setup MDP.
6  mdp = GridWorldMDP(width=4, height=3, init_loc=(1, 1), goal_locs=[[4, 3]])
7
8  # Make agents.
9  ql_agent = QLearningAgent(actions=mdp.get_actions())
10 rand_agent = RandomAgent(actions=mdp.get_actions())
11
12 # Run experiment and make plot.
13 run_agents_on_mdp([ql_agent, rand_agent], mdp, instances=5, episodes=50,
    steps=10)

```

---

Figure 2: Example code for running a basic experiment. First, define a grid-world MDP (line 6), then make our agents (line 9-10), and then run the experiment (line 13). Running the above will generate the plot shown in Figure 4.

## 12 1 Introduction

13 Reinforcement learning (RL) has recently soared in popularity due in large part to recent success  
14 in challenging domains, including learning to play Atari games from image input [27], beating the  
15 world champion in Go [32], and robotic control from high dimensional sensors [21]. In concert with  
16 the field’s growth, experiments have become more complex, leading to new challenges for empir-  
17 ical evaluation of RL methods. Recent work by Henderson et al. [16] highlighted many of the is-  
18 sues involved with handling this new complexity, raising concerns about emerging RL experimental  
19 practices. Additionally, Python has become a prominent programming language used by machine-  
20 learning researchers due to the availability of powerful deep learning libraries like PyTorch [29] and  
21 tensorflow [1], along with scipy [19] and numpy [28].

22 To accommodate this growth, there is a need for a simple, lightweight library that supports quick  
23 execution and analysis of RL experiments in Python. Certainly, many libraries already fulfill this  
24 need for many uses cases—as will be discussed in Section 2, many effective RL libraries for Python  
25 already exist. However, the design philosophy and ultimate end user of these packages is distinct  
26 from those targeted by `simple_rl`: those users who seek to quickly run simple experiments, look at  
27 a plot that summarizes results, and allow for the quick sharing and reproduction of these findings.

28 The core design principle of `simple_rl` is that of *simplicity*, per its name. The library is stripped  
29 down to the bare necessities required to run basic RL experiments. The focus of the library is on  
30 traditional, tabular domains, though it does have the capacity to cooperate with high-dimensional  
31 environments like those offered by the OpenAI Gym [6]. The assumed objective of a practitioner  
32 using the library is to define (1) an RL agent (or collection of agents), (2) an environment (an  
33 MDP, POMDP, or similar Markov model), (3) let the agent(s) interact with the environment, and  
34 (4) view and analyze the results of this interaction. This basic pipeline serves as the “end-game” of  
35 `simple_rl`, and dictates much of the design and its core features. A block diagram of this process  
36 is presented in Figure 1: run an experiment, see the results, and reproduce these results according  
37 to an auto-generated JSON file logging the experimental details. The actual code of the experiment  
38 run is shown in Figure 2: in around five lines, we define a  $Q$ -Learning instance, a random actor, and  
39 a simple grid-world domain, and let these agents interact with the environment for a set number of  
40 instances. As mentioned, running this code produces both a JSON file tracking the experiment that  
41 can be used (or shared) to run the same experiment again, and regenerate the plot seen in Figure 4a.

## 42 2 Relation To Other Libraries

43 Many excellent libraries already exist in Python for carrying out RL experiments. What separates  
44 `simple_rl`? As the name suggests, its distinguishing feature is its emphasis on simplicity, which  
45 also brings a shortage of certain features. We here describe the objectives of other RL libraries in

46 Python, and briefly cover what some have implemented in case those are a better fit for the needs of  
47 different programmers.

## 48 **2.1** RLPy

49 RLPy offers a well documented, expansive library for RL and planning experiments in Python 2 [15].  
50 The library includes a similar overall structure to that of `simple_rl`: the core entities are agents,  
51 environments, experiments, policies, and representations. The main focus of RLPy is on value-  
52 function approximation, but the library also offers several MDP solvers in the form of the usual  
53 dynamic programming algorithms like value iteration [4] and policy iteration [18]. Notably, the  
54 library also includes a large number of canonical RL tasks, including Mountain Car, Acrobot, Puddle  
55 World, Swimmer, and Cart Pole.

56 Get it here: <https://github.com/rlpy/rlpy>

## 57 **2.2** mushroom

58 Mushroom is a new library aimed at simplifying RL experimentation with OpenAI gym and tensor-  
59 flow, but also offers support for traditional tabular experiments [13]. Mushroom offers implemen-  
60 tations of many recent Deep RL algorithms, including DQN [27], Stochastic Actor-Critic [12], and  
61 a template for Policy Gradient algorithms. All of its neural network code is based on tensorflow.  
62 Additionally, Mushroom comes with noteworthy RL tasks like Mountain Car, Inverted Pendulum,  
63 and a classic Linear-Quadratic Regulator control task.

64 Get it here: <https://github.com/AIRLab-POLIMI/mushroom>

## 65 **2.3** PyBrain

66 PyBrain is an established, expansive, general purpose library for machine learning in Python [30],  
67 but also offers infrastructure for conducting RL experiments with a similar focus to RLPy. The  
68 library includes a number of the standard environments and agents, with a large number of model-  
69 free algorithms.

70 Get it here: <http://www.pybrain.org/>

## 71 **2.4** keras-rl

72 `keras-rl` provides integration between Keras [9] and many popular Deep RL algorithms.  
73 `keras-rl` offers an expansive list of implemented Deep RL algorithms in one place, including:  
74 DQN, Double DQN [37], Deep Deterministic Policy Gradient [23], and Dueling DQN [38]. For  
75 those that use Keras for deep learning and mostly want to focus on deep RL, `keras-rl` library is a  
76 great choice.

77 Get it here: <https://github.com/keras-rl/keras-rl>

## 78 **2.5** python-rl

79 `python-rl` [11] provides integration with the classic language-agnostic framework RL-Glue [36].  
80 The main goal of this library is to bring RL-Glue up to date with a few somewhat more recent  
81 features, agents, and environments in common RL experiments.

82 Get it here: <https://github.com/amarack/python-rl>

## 83 **2.6** reinforcement-learning

84 `reinforcement-learning` offers an excellent resource for RL education—it is designed to  
85 be paired with David Silver’s online RL course<sup>2</sup> [5]. The library contains many central al-  
86 gorithms, including value iteration, policy iteration, Q-Learning [39], SARSA [33], and Pol-

---

<sup>2</sup><https://www.youtube.com/watch?v=2pWv7G0vuf0>

87 icy Gradient [40, 35]. Programmers planning to go through David Silver’s course may find the  
88 reinforcement-learning library the most suitable package.

89 Get it here: <https://github.com/dennybritz/reinforcement-learning>

## 90 2.7 dopamine

91 dopamine is a recently released library [3] offering many of the most recent deep RL algorithms  
92 including Rainbow [17], Prioritized Experience Replay [31], and Distributional RL [2], with an eye  
93 for reproducibility in the ALE based on the suggestions given by [25]. dopamine offers a lot for  
94 people whose main agenda is to run experiments in the ALE or perform new research in deep RL.

95 Get it here: <https://github.com/google/dopamine>

96

.....

97 To summarize: Many great packages are already out there. The main differentiating features of  
98 simple\_rl are (1) quick generation of plots, (2) focus on reproducibility, and (3) emphasis on  
99 simplicity, both in terms of algorithmic development and its attachment to classical RL problems  
100 (like grid worlds).

## 101 3 Overview of Features

102 We begin by unpacking the example in Figure 2 to showcase the main design philosophy of  
103 simple\_rl.

### 104 3.1 The Core: Agents and MDPs

105 The library primarily consists of *agents* and *environments* (called “tasks” in the library).

106 Agents, by default, are all subclasses of the abstract class, Agent, which is only responsible for  
107 a method `act(self, state, reward) → action`. A list of agents, planning algorithms, and  
108 tasks currently implemented is presented in Table 1.

109 Tasks, for the most part, all inherit from the abstract MDP class, MDP. The core of an MDP is  
110 its transition function and reward function, captured in the abstract class by class-wide variables,  
111 `transition_func` and `reward_func`:

`transition_func(state, action) → state,` (1)

`reward_func(state, action) → reward.` (2)

112 When defining an MDP instance, you must pass in functions of  $T$  and  $R$  that *output* a state and re-  
113 ward, respectively. In this way, no MDP is ever responsible for enumerating either  $\mathcal{S}$  or  $\mathcal{A}$  explicitly,  
114 thereby allowing for (1) simple specification of these two functions, and (2) efficient implementation  
115 of high-dimensional domains—we need only represent and store the states that are visited during  
116 experimentation.

117 Naturally, MDP subclasses have a variety of arguments—in the earlier grid-world example, we saw  
118 the `GridWorldMDP` class take as input the dimensions of the grid, a starting location, and a list of  
119 goal locations. Such inputs are typical to MDP classes in simple\_rl.

#### 120 3.1.1 Running Simple Experiments

121 Defining an agent and an MDP is almost all that is needed to run an experiment. The final component  
122 required is an experiment function from the `run_experiments.py` file. This file contains a number  
123 of different experiment types that are catered to the different environment types (POMDPs, Markov  
124 Games, and so on). For now, let us focus on `run_agents_on_mdp` function, which is the most  
125 canonical. As per the example in Figure 2, this function takes as input at minimum a list of agents  
126 and an MDP instance. A user can also specify experimental parameters like `instances`, `episodes`,  
127 and `steps`, which indicate the following:

- 128 • `instances`: The number of times to repeat the entire experiment (will be used to form  
129 95% confidence intervals for all experiments conducted).

<i>RL Agents</i>	Q-Learning, RMax, DelayedQ, DoubleQ, Random, Fixed Linear Q-Learning, DQN, LinUCB.
<i>Planning Algorithms</i>	Value Iteration, Bounded RTDP, MCTS
<i>MDPs</i>	Chain, Grid World, Randomized Graph, Open AI Gym Combo Lock, Puddle, Hanoi, Bandit
<i>OOMDPs</i>	Taxi, Trench, Cleanup
<i>POMDPs</i>	Maze
<i>Markov Games</i>	Grid Games, Rock Paper Scissors, Prisoner’s Dilemma, Gather

Table 1: An overview of Agents and MDPs in `simple_rl`.

```

1 from simple_rl.tasks import GymMDP
2 from simple_rl.agents import RandomAgent, LinearQAgent
3 from simple_rl.run_experiments import run_agents_on_mdp
4
5 # Gym MDP
6 gym_mdp = GymMDP(env_name='CartPole-v1', render=True)
7 num_feats = gym_mdp.get_num_state_feats()
8
9 # Setup agents and run.
10 rand_agent = RandomAgent(gym_mdp.get_actions())
11 lin_q_agent = LinearQAgent(gym_mdp.get_actions(), num_feats, rbf=True)
12 agents = [lin_q_agent, rand_agent]
13
14 # Run.
15 run_agents_on_mdp(agents, gym_mdp, instances=5, episodes=5000, steps=200)

```

Figure 3: Running experiments in the OpenAI Gym.

130 • *episodes*: The number of *episodes* per instance. An episode will consist of *steps* number  
131 of steps, after which the agent is reset to the start state (but gets to remember what it has  
132 learned so far).

133 • *steps*: The number of steps per episode.

134 The plotting is set up to plot all of the above appropriately. For instance, if a user sets `episodes=1`  
135 but `steps=50`, then the library produces a step-wise plot (that is, the x-axis is steps, not episodes).

136 Running the function `run_agents_on_mdp` will create a JSON file detailing all of the components  
137 of the experiment needed to rerun it. Then, it will create a folder locally, “results”, store each  
138 agent’s stream of received rewards, and print out the status of the experiment to console. When  
139 the experiment concludes, a learning curve with 95% confidence intervals will be generated (via  
140 `simple_rl/utils/chart_utils.py` and opened. The JSON file lets users of the library recon-  
141 struct and rerun the original experiment using another function from the `run_experiments.py`  
142 script. In this way, the JSON file is effectively a certificate that this plot can be reproduced if the  
143 same experiment were run again. We provide more detail on this feature in Section 3.2.

144 We can also run a similar experiment in the OpenAI Gym (Figure 3).

145 As can be seen in Figure 3, the structure of the experiment is identical. Since we define a `GymMDP`,  
146 we pass as input the name of the environment we’d like to produce: In this case, we’re running ex-  
147 periments in `CartPole-v1`, but any of the usual Gym environment names will work. We can also pass  
148 in the `render` boolean flag, indicating whether or not we’d like to visualize the learning process. Al-  
149 ternatively, we can pass in the `render_every_n_episodes` flag (along with `render=True`), which  
150 will only render the agent’s learning process every  $N$  episodes.

151 On longer experiments, we may want additional feedback about the learning process. For this purpose, the `run_agents_on_mdp` function also takes as input a Boolean flag `verbose`, which, if true, will provide detailed episode-by-episode tracking of the progress of the experiment to the console. 152  
153 There are a number of other ways to run experiments, but these examples capture the core experimental cycle. 154  
155

156 **Other Environment Types** The library offers support for other types of environments beyond typical MDPs, including classes for Object-Oriented MDPs or OOMDPs [14],  $k$ -Armed Bandits [8], 157 Partially Observable MDPs or POMDPs [20], a probability distribution over MDPs for lifelong learning [7], and Markov Games [24]. Aspects of these classes are handled slightly differently to 158 accommodate the different kinds of decision-making problems they capture, but the interface to run 159 experiments with each type is nearly identical. Examples for how to run experiments with each type 160 of environment are included in the examples directory in the repository along with a test script that 161 ensures each example can run on a given machine. Running experiments with these other environment 162 types is the same as the pipeline so far described: a function in the `run_experiments.py` 163 script will handle all of the interactions between agent(s) and environment and produce a plot when 164 the experiment finishes. Notably, the reproducibility feature is not yet fully developed for all environment 165 types. This is a major direction for future development of the library. 166  
167

### 168 3.2 Reproducibility

169 Due to its simplicity, the library is naturally suited for reproducing results from previously run experiments. As mentioned, *every* experiment that is conducted using the library will create a directory 170 with the experiment name containing a JSON file “`full_experiment_data.json`” that enumerates every parameter, agent, MDP, and type needed to launch the exact same experiment another 171 time. The idea is that these files can be shared across users of the library—if a user gives someone else this file (and the necessary agents and environments), it is a contract that they can rerun *exactly* 172 the same experiment just run using `simple_rl`. 173  
174

175 Using one of these experiment files, the function `reproduce_from_exp_file(exp_name)`, will 176 read the experiment file, reconstruct all the necessary components, rerun the entire experiment, and 177 remake the plot. Thus, providing one of these JSON files is to be interpreted as a certificate that this 178 experiment is guaranteed to produce similar results. 179

180 As an example, consider again the code from Figure 2. Running this code will create: (1) 181 the “`results`” directory, (2) the “`gridworld_h-3_w-4`” directory within `results`, and (3) the 182 “`full_experiment_data.json`” file, which contains *all* necessary parameters to rerun the experiment. 183

184 Suppose someone provided the directory `gridworld_h-3_w-4` containing the experiment file for 185 the above grid-world experiment. Then, we could run the following code:

```
186  
187 1 from simple_rl.run_experiments import reproduce_from_exp_file  
188 2  
189 3 reproduce_from_exp_file("gridworld_h-3_w-4")
```

191 Which will automatically generate the plot in Figure 4b.

192 To ensure reproducibility of new subclasses or other bells and whistles attached to the library, any 193 agent or MDP must implement the “`get_parameters(self)`” method that returns a dictionary 194 containing all relevant parameters for the instance to be reconstructed. For example, consider the 195 `QLearningAgent` class in Figure 5.

196 Any introduced subclass that wants to play along well with the reproduction infrastructure in 197 `simple_rl` must have such a method.

198 We stipulate that this is a lightweight means of ensuring reproduction for three reasons: 1) it is 199 entirely obfuscated from the programmer, as all tracking of experimental parameters is done automatically, 2) a single, universally formatted document (JSON) contains all the information needed to 200 guarantee reproduction of results (along with a copy of the library itself, and any new agents/MDPs), 201 and 3) the library is simple enough that most experiments consist of only a small number of moving 202

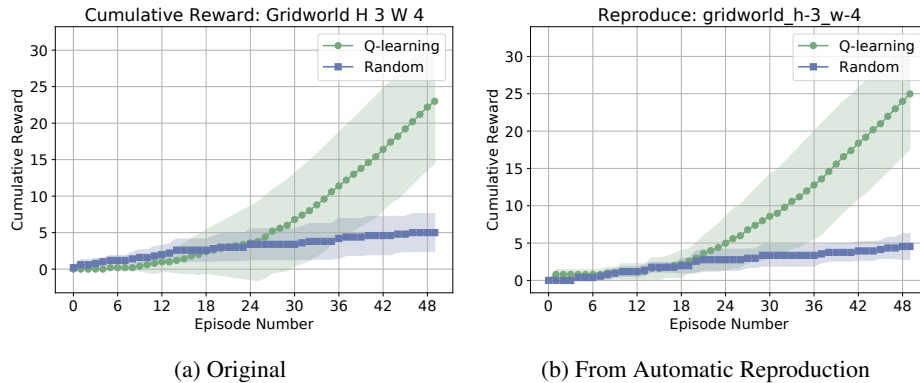


Figure 4: Original results (left) and results generated by reproducing the experiment (right).

---

```

1  def get_parameters(self):
2      """
3      Returns:
4          (dict) key=param_name (str) --> val=param_val (object).
5      """
6      param_dict = defaultdict(int)
7
8      param_dict["alpha"] = self.alpha
9      param_dict["gamma"] = self.gamma
10     param_dict["epsilon"] = self.epsilon_init
11     param_dict["anneal"] = self.anneal
12     param_dict["explore"] = self.explore
13
14     return param_dict

```

---

Figure 5: The `get_parameters` method of `QLearningAgent`.

203 parts. The feature to reproduce from a JSON does not yet fully support all environment types, but it  
 204 is an active area of development for the library.

205 To recap, the introduced components define the essence of the library:

- 206 • Center everything around *agents*, *MDPs*, and interactions thereof.
- 207 • Completely obscure the complexity of plotting and experiment tracking from the program-  
 208 mer, while making it simple to plot and reproduce results if needed.
- 209 • Simplicity above all else.
- 210 • Treat things *generatively*—namely, MDPs transition models and reward functions are best  
 211 implemented as functions that return a state or reward, rather than enumerate all state-  
 212 actions pairs.

### 213 3.3 Utilities

214 In addition to the core experimental pipeline described above, the library is well stocked with other  
 215 utilities useful for RL and planning.

216 **Plotting** As is shown by Figure 1, plotting is tightly coupled with running experiments.  
 217 Each experiment type is connected with the same plotting script, stored in the library in  
 218 `utils/chart_utils.py`. The basic plot shows some measure of time along the x-axis (either in  
 219 episodes run or steps taken), with cumulative reward shown in the y-axis for each given algorithm.  
 220 While this plot is the default learning curve generated, the experimental pipeline gives the end pro-  
 221 grammer control over the type of plot generated. First, the flag `cumulative_plot` for all of the core

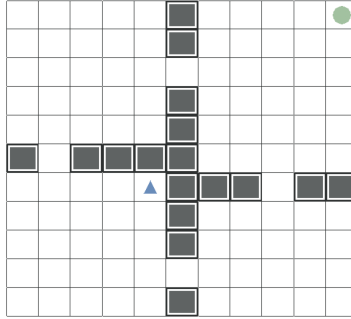


Figure 6: Example visual generated by the library

222 experiment functions is set to True by default (as in `run_agents_on_mdp`, `run_agents_lifelong`).  
 223 Thus, if we simply run the experiment with this flag set to False, we'll produce an average reward  
 224 plot instead. Second, the default y-axis is cumulative *reward*—sometimes, though, we'd like to  
 225 measure the *discounted* reward acquired by the agent. To do so, we set the `track_disc_reward`  
 226 flag of any of the core experimental functions to True. There are also mechanisms for plotting the  
 227 wall-clock time taken by each agent, and plotting the percentage of successful runs of each agent,  
 228 where success is defined according to a user defined function on the reward stream received by the  
 229 agent. For more details on plotting, see the `chart_utils.py` script.

230 **Visuals** The library offers bare bones visuals for the grid world domains using `pygame`<sup>3</sup>. An  
 231 example is presented in Figure 6; in this case, the learning process is visualized while the experiment  
 232 runs. The library also supports visualizing policies and value functions, so long as an MDP comes  
 233 along with a `draw_state` method, and an interactive mode where the user can control the agent via  
 234 keyboard input. However, visuals are very much an underdeveloped aspect of the library. A major  
 235 point of future development is to equip `simple_rl` with a comprehensive suite of visualization and  
 236 analysis tools.

237 **Abstraction** A core approach to RL involves forming *abstractions*, either of state [22] or ac-  
 238 tion [34]. `simple_rl` contains support for planning or learning with either state aggregation func-  
 239 tions, which compress a given MDP's state space into a smaller one, and options, which encode long  
 240 horizon sequences of actions, useful for targeted exploration and efficient planning.

241 **Planning** The library includes several default planning algorithms such as Value Iteration, Monte  
 242 Carlo Tree Search [10], and Bounded Real Time Dynamic Programming [26]. Planners can be used  
 243 to compute the value function, the optimal (or near-optimal) policy, or enumerate a state-action  
 244 space (see `planning_example.py` in the repository).

## 245 4 Conclusion

246 `simple_rl` offers a lightweight suite of tools for conducting RL experiments in Python 2 and 3.  
 247 Its design philosophy focuses on obfuscating complexity from the end user, including the tracking  
 248 of experimental details, generation of plots, and construction of agents and MDPs. This leads to a  
 249 package that is relatively light in features but comes with an ease of use that lets only a few lines  
 250 of code generate learning curves that are guaranteed to be reproducible. The library is available  
 251 on the Python package index, and thus can be installed with the usual `pip install simple_rl`.  
 252 In progress documentation is available as well.<sup>4</sup> Many features are currently under development:  
 253 the most important near term goal is to expand the suite of reproducibility tools to account for  
 254 more variety across different operating systems and other variables that might impact experiments.  
 255 Additionally, the library lacks a suite of basic deep RL algorithms for use in experimentation, a  
 256 general interface for visualizing MDPs (and other environments), and a more expansive collection  
 257 of tasks, RL algorithms, and planning algorithms.

<sup>3</sup><https://pygame.org>

<sup>4</sup>[https://david-abel.github.io/simple\\_rl/docs/index.html](https://david-abel.github.io/simple_rl/docs/index.html)



## References

- 258
- 259 [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu  
260 Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: a system for  
261 large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- 262 [2] Marc G. Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on rein-  
263 forcement learning. In *Proceedings of the International Conference on Machine Learning*,  
264 volume 70, pages 449–458, 2017.
- 265 [3] Marc G. Bellemare, Pablo Samuel Castro, Carles Gelada, Saurabh Kumar, and Subhodeep  
266 Moitra. dopamine, 2018.
- 267 [4] Richard Bellman. A Markovian decision process. *Journal of Mathematics and Mechanics*, 6  
268 (5):679–684, 1957.
- 269 [5] Denny Britz. reinforcement-learning. [github.com/dennybritz/  
270 reinforcement-learning](https://github.com/dennybritz/reinforcement-learning), 2018.
- 271 [6] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang,  
272 and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 273 [7] Emma Brunskill and Lihong Li. Pac-inspired option discovery in lifelong reinforcement learn-  
274 ing. In *International Conference on Machine Learning*, pages 316–324, 2014.
- 275 [8] Robert R Bush and Frederick Mosteller. A stochastic model with applications to learning. *The  
276 Annals of Mathematical Statistics*, pages 559–585, 1953.
- 277 [9] François Chollet et al. Keras. <https://keras.io>, 2015.
- 278 [10] Rémi Coulom. Efficient selectivity and backup operators in Monte-Carlo tree search. In *Pro-  
279 ceedings of the International Conference on Computers and Games*, pages 72–83. Springer,  
280 2006.
- 281 [11] Will Dabney and Pierre-Luc Bacon. python-rl. [https://github.com/amarack/  
282 python-rl](https://github.com/amarack/python-rl), 2013.
- 283 [12] Thomas Degris, Patrick M Pilarski, and Richard S Sutton. Model-free reinforcement learning  
284 with continuous action in practice. In *American Control Conference (ACC), 2012*, pages 2177–  
285 2182. IEEE, 2012.
- 286 [13] Carlo D’Eramo and Davide Tateo. mushroom. [https://github.com/AIRLab-POLIMI/  
287 mushroom](https://github.com/AIRLab-POLIMI/mushroom), 2018.
- 288 [14] Carlos Diuk, Andre Cohen, and Michael L Littman. An object-oriented representation for  
289 efficient reinforcement learning. In *Proceedings of the International Conference on Machine  
290 Learning*, pages 240–247. ACM, 2008.
- 291 [15] Alborz Geramifard, Robert H Klein, Christoph Dann, William Dabney, and Jonathan P How.  
292 RLPy: The reinforcement learning library for education and research. [http://acl.mit.  
293 edu/RLPy](http://acl.mit.edu/RLPy), 2013.
- 294 [16] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David  
295 Meger. Deep reinforcement learning that matters. 2018.
- 296 [17] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dab-  
297 ney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining im-  
298 provements in deep reinforcement learning. *Proceedings of the AAAI Conference on Artificial  
299 Intelligence*, 2018.
- 300 [18] Ronald A Howard. Dynamic programming and Markov processes. 1960.
- 301 [19] Eric Jones, Travis Oliphant, and Pearu Peterson. Scipy: Open source scientific tools for python.  
302 2014.

- 303 [20] Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in  
304 partially observable stochastic domains. *Artificial Intelligence*, 101(1-2):99–134, 1998.
- 305 [21] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep  
306 visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- 307 [22] Lihong Li, Thomas J Walsh, and Michael L Littman. Towards a unified theory of state abstrac-  
308 tion for MDPs. In *ISAIM*, 2006.
- 309 [23] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval  
310 Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning.  
311 2016.
- 312 [24] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In  
313 *Machine Learning*, pages 157–163. Elsevier, 1994.
- 314 [25] Marlos C Machado, Marc G Bellemare, Erik Talvitie, Joel Veness, Matthew Hausknecht, and  
315 Michael Bowling. Revisiting the arcade learning environment: Evaluation protocols and open  
316 problems for general agents. *arXiv preprint arXiv:1709.06009*, 2017.
- 317 [26] H Brendan McMahan, Maxim Likhachev, and Geoffrey J Gordon. Bounded real-time dynamic  
318 programming: RTDP with monotone upper bounds and performance guarantees. In *Proceed-*  
319 *ings of the International Conference on Machine Learning*, pages 569–576. ACM, 2005.
- 320 [27] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G  
321 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.  
322 Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- 323 [28] Travis Oliphant. *Guide to NumPy*. 01 2006.
- 324 [29] Adam Paszke, Sam Gross, Soumith Chintala, and Gregory Chanan. Pytorch: Tensors and  
325 dynamic neural networks in python with strong gpu acceleration, 2017.
- 326 [30] Tom Schaul, Justin Bayer, Daan Wierstra, Yi Sun, Martin Felder, Frank Sehnke, Thomas  
327 Rückstieß, and Jürgen Schmidhuber. PyBrain. *Journal of Machine Learning Research*, 2010.
- 328 [31] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay.  
329 *arXiv preprint arXiv:1511.05952*, 2015.
- 330 [32] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van  
331 Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot,  
332 et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529  
333 (7587):484, 2016.
- 334 [33] Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine*  
335 *Learning*, 3(1):9–44, 1988.
- 336 [34] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A frame-  
337 work for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–  
338 211, 1999.
- 339 [35] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradi-  
340 ent methods for reinforcement learning with function approximation. In *Advances in Neural*  
341 *Information Processing Systems*, pages 1057–1063, 2000.
- 342 [36] Brian Tanner and Adam White. RI-glue: Language-independent software for reinforcement-  
343 learning experiments. *Journal of Machine Learning Research*, 10(Sep):2133–2136, 2009.
- 344 [37] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double  
345 q-learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2016.
- 346 [38] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Van Hasselt, Marc Lanctot, and Nando De Fre-  
347 itas. Dueling network architectures for deep reinforcement learning. *Proceedings of the Inter-*  
348 *national Conference on Machine Learning*, 2016.

- 349 [39] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8(3-4):279–292,  
350 1992.
- 351 [40] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist rein-  
352 forcement learning. *Machine Learning*, 8(3-4):229–256, 1992.