

---

# Prioritizing Samples in Reinforcement Learning with Reducible Loss

---

Anonymous Author(s)

Affiliation

Address

email

## 1 A Implementation Details

2 We build our experiments on top of existing implementations of SAC, DQN and Rainbow. For the  
3 DeepMind Control Suite experiments, we modify Yarats and Kostrikov [2020], adding a prioritized  
4 replay buffer and the ReLo version. We use an open source implementation of Rainbow<sup>1</sup> for the  
5 Arcade Learning Environment and the DQN implementation from the MinAtar authors Young and  
6 Tian [2019]. Aside from the collected frames and number of seeds, we have not modified any of the  
7 hyper-parameters from these original implementations. The hyper-parameters as well as hardware  
8 and software used are given in Table 1.

Table 1: Hyper-Parameters of all experiments

Environments	Algorithm	Algorithm Parameters	Hardware & Software
ALE	Rainbow	Frames = $2 \times 10^6$ seeds = 5	Hardware- CPU: 6 Intel Gold 6148 Skylake GPU: 1 NVidia V100 RAM: 32 GB
		Remaining hyper-parameters same as Hessel et al. [2017]	Software- Pytorch: 1.10.0 Python: 3.8
DeepMind Control Suite	SAC	Frames = $1 \times 10^6$ seeds = 5	Hardware- CPU: 6 Intel Gold 6148 Skylake GPU: 1 NVidia V100 RAM: 32 GB
		Remaining hyper-parameters same as Haarnoja et al. [2018]	Software- Pytorch: 1.10.0 Python: 3.8
MinAtar	DQN	Frames = $5 \times 10^6$ seeds = 5	Hardware- CPU: 6 Intel Gold 6148 Skylake GPU: 1 NVidia V100 RAM: 32 GB
		Remaining hyper-parameters same as Mnih et al. [2015]	Software- Pytorch: 1.10.0 Python: 3.8

---

<sup>1</sup><https://github.com/Kaixhin/Rainbow>

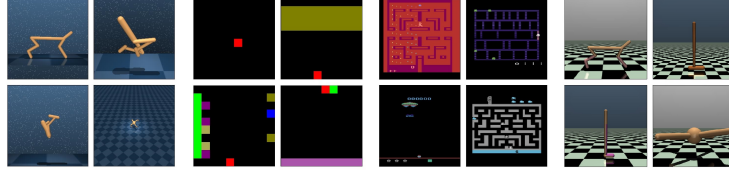


Figure 1: Visualization of a few environments from each benchmark. Left to right: DeepMind Control Suite, MinAtar, Arcade Learning Environment

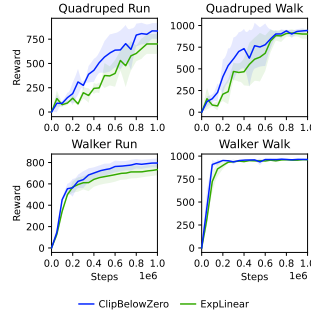


Figure 2: Comparison of different mapping functions from ReLo to  $p_i$  on a subset of environments from the DMC benchmark. Performance is evaluated over 3 seeds.

## 9 B Mapping Functions for ReLo

10 Prioritized experience replay buffers expect the priorities assigned to data points to be non-negative.  
 11 While the MSE version of the TD error used in PER satisfies this constraint, ReLo does not. Therefore,  
 12 there must be a non-negative, monotonically increasing mapping from ReLo to  $p_i$ . In our main  
 13 experiments, we clipped negative ReLo values to zero. Another mapping we tried was to set  
 14  $p_i = e^{\text{ReLo}}$ , in which case the probability of sampling a data point  $P_i$ , from Eq. 6, corresponds to the  
 15 softmax over ReLo scores. However, for this choice the priority would explode if the ReLo crossed  
 16 values above 40 which happened occasionally during the initial stages of learning in Rainbow. The  
 17 second mapping function candidate was exponential when ReLo is negative and linear otherwise,  
 18 that is,

$$f_{\text{ExpLinear}} = \begin{cases} e^{\text{ReLo}} & \text{if ReLo} < 0 \\ \text{ReLo} + 1 & \text{otherwise} \end{cases}$$

19 The linear part is shifted so that the mapping is smooth around  $\text{ReLo} = 0$ . As shown in Fig. 2,  
 20  $f_{\text{ExpLinear}}$  performs worse compared to just clipping negative ReLo values to zero. When the ReLo  
 21 values during training are analysed, we observe that the average of ReLo values (before the mapping)  
 22 tends to be positive, so clipping does not lead to a large loss in information.

## 23 C Extended Related Work

### 24 C.1 Prioritization Schemes

25 Prioritization strategies have been leading to important improvements in sample efficiency. Sinha  
 26 et al. [2020] proposes an approach that re-weights experiences based on their likelihood under the  
 27 stationary distribution of the current policy in order to ensure small approximation errors on the value  
 28 function of recurring seen states. Lahire et al. [2021] introduces the Large Batch Experience Replay  
 29 (LaBER) to overcome the issue of the outdated priorities of PER and its hyperparameter sensitivity by  
 30 employing an importance sampling view for estimating gradients. LaBER first samples a large batch  
 31 from the replay buffer then computes the gradient norms and finally down-samples it to a mini-batch  
 32 of smaller size according to a priority.

33 Kumar et al. [2020] presents Distribution Correction (DisCor), a form of corrective feedback to  
 34 make learning dynamics more steady. DisCor computes the optimal distribution and performs a

35 weighted Bellman update to re-weight data distribution in the replay buffer. Inspired by DisCor,  
 36 Regret Minimization Experience Replay (ReMERN) Liu et al. [2021] estimates the suboptimality  
 37 of the  $Q$  value with an error network. Yet, Hong et al. [2022] uses Topological Experience Replay  
 38 (TER) to organize the experience of agents into a graph that tracks the dependency between  $Q$ -value  
 39 of states.

40 While PER was initially proposed as an addition to DQN-style agents, Hou et al. [2017] have shown  
 41 that PER can be a useful strategy for improving performance in Deep Deterministic Policy Gradients  
 42 (DDPG) Lillicrap et al. [2016]. Another recent strategy to improve sample efficiency was to introduce  
 43 losses from the transition dynamics along with the TD error as the priority Oh et al. [2022]. Although  
 44 this has shown improvements, it involves additional computational complexity since it also requires  
 45 learning a reward predictor and transition predictor for the environment. Our proposal does not  
 46 require training additional networks and hence is similar in computational complexity to PER. This  
 47 makes it very simple to integrate into any existing algorithm. Wang and Ross [2019] propose an  
 48 algorithm to dynamically reduce the replay buffer size during training of SAC so that the agent  
 49 prioritizes recent experience while also ensuring that updates performed using newer data are not  
 50 overwritten by updates from older data. However, they do not distinguish between points based on  
 51 learn-ability and only assume that newer data is more useful for the agent to learn.

## 52 C.2 Off-Policy Algorithms

53 Off-policy algorithms are those that can learn a policy by learning from data not generated from  
 54 the current policy. This improves sample efficiency by reusing data collected by old versions of  
 55 the policy. This is in contrast to on-policy algorithms such as PPO Schulman et al. [2017], which  
 56 after collecting a batch of data and training on it, discard those samples and start data collection  
 57 from scratch. Recent state-of-the-art off-policy algorithms for continuous control include Soft Actor  
 58 Critic (SAC) Haarnoja et al. [2018] and Twin Delayed DDPG (TD3) Fujimoto et al. [2018]. SAC  
 59 learns two  $Q$  networks together and uses the minimum of the  $Q$  values generated by these networks  
 60 for the Bellman update equation to avoid over estimation bias. The  $Q$  target update also includes a  
 61 term to maximize the entropy of the policy to encourage exploration, a formulation that comes from  
 62 Maximum Entropy RL Ziebart et al. [2008]. TD3 is a successor to DDPG Lillicrap et al. [2016]  
 63 which addresses the overestimation bias present in DDPG in a similar fashion to SAC, by learning  
 64 two  $Q$  networks in parallel, which explains the “twin” in the name. It learns an actor network  $\mu$   
 65 following Eq. 4 to compute the maximum over  $Q$  values. TD3 proposes that the actor networks be  
 66 updated at a less frequent interval than the  $Q$  networks, which gives rise to the “delayed” name. In  
 67 discrete control, Rainbow Hessel et al. [2017] combines several previous improvements over DQN,  
 68 such as Double DQN van Hasselt et al. [2016], PER Schaul et al. [2016], Dueling DQN Wang et al.  
 69 [2016], Distributional RL Bellemare et al. [2017] and Noisy Nets Fortunato et al. [2018].

## 70 D DeepMind Control Suite

71 We choose 9 environments from the DeepMind Control Suite Tassa et al. [2018] for testing the  
 72 performance of ReLo on continuous control tasks. Each agent was trained on proprioceptive inputs  
 73 from the environment for 1M frames with an action repeat of 1. The training curves for the baselines  
 74 and ReLo are given in Fig. 3.

## 75 E OpenAI Gym Environments

76 We evaluate agents for 1M timesteps on each environment and similar to DM Control, they are  
 77 trained using proprioceptive inputs from the environment. The hyperparameters for this benchmark  
 78 are shared with those used for the DM Control Suite experiments.

## 79 F MinAtar

80 We evaluate the baselines against all 5 environments in the MinAtar suite Young and Tian [2019].  
 81 A visualization of a few environments from the suite is presented in Fig. 1. Each agent receives  
 82 the visual observations from the environment and is trained for 5M frames following the evaluation  
 83 methodology outlined in Young and Tian [2019]. The training curves are given in Fig. 5.

Table 2: Comparison of PER and ReLo on the DMC benchmark

	BASELINE	PER	LABER	ReLo
CHEETAH RUN	761.89 $\pm$ 112.38	<b>831.87 <math>\pm</math> 38.90</b>	579.80 $\pm$ 60.61	660.29 $\pm$ 141.22
FINGER SPIN	966.68 $\pm$ 29.40	975.40 $\pm$ 6.75	884.46 $\pm$ 20.23	<b>978.78 <math>\pm</math> 14.46</b>
HOPPER HOP	<b>264.75 <math>\pm</math> 37.90</b>	217.35 $\pm$ 113.79	90.85 $\pm$ 57.76	247.81 $\pm$ 51.01
QUADRUPED RUN	612.67 $\pm$ 143.90	496.42 $\pm$ 216.01	544.80 $\pm$ 41.84	<b>833.92 <math>\pm</math> 81.05</b>
QUADRUPED WALK	831.92 $\pm$ 74.34	766.30 $\pm$ 200.86	716.61 $\pm$ 270.36	<b>942.64 <math>\pm</math> 9.75</b>
REACHER EASY	<b>983.06 <math>\pm</math> 2.70</b>	981.58 $\pm$ 6.33	947.20 $\pm$ 14.46	979.08 $\pm$ 11.02
REACHER HARD	955.08 $\pm$ 38.52	935.08 $\pm$ 47.94	951.08 $\pm$ 6.70	<b>956.80 <math>\pm</math> 38.73</b>
WALKER RUN	759.13 $\pm$ 23.91	755.49 $\pm$ 64.35	551.81 $\pm$ 58.41	<b>795.14 <math>\pm</math> 42.52</b>
WALKER WALK	943.67 $\pm$ 30.28	957.38 $\pm$ 8.24	863.89 $\pm$ 112.60	<b>963.28 <math>\pm</math> 5.03</b>

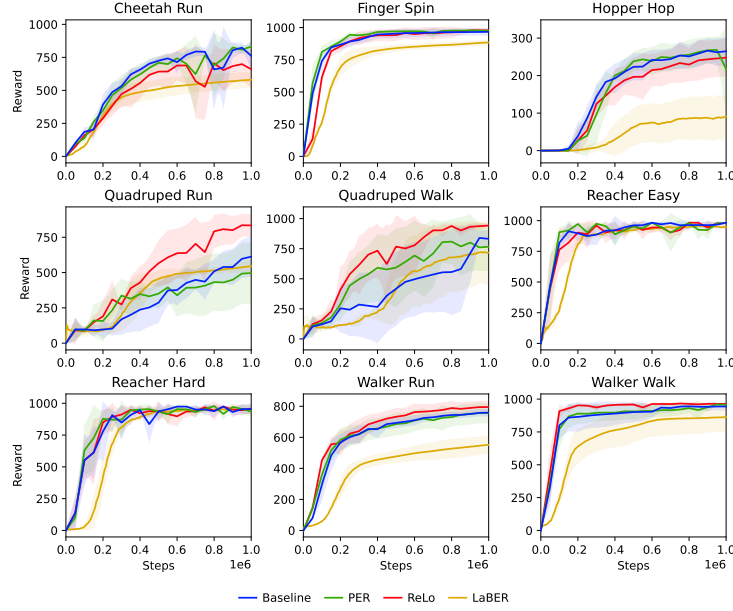


Figure 3: Training curves of environments from the DeepMind Control Suite. Performance is evaluated for 10 episodes over 5 random seeds.

Table 3: Comparison of PER and ReLo on OpenAI Gym environments

	BASELINE	PER	ReLo
GYM HALFCHEETAH	9579.60 $\pm$ 1331.00	9549.42 $\pm$ 917.92	<b>11590.63 <math>\pm</math> 670.36</b>
GYM HOPPER	<b>2795.18 <math>\pm</math> 659.19</b>	2175.09 $\pm$ 456.11	2527.93 $\pm$ 648.61
GYM WALKER	2854.20 $\pm$ 623.69	735.16 $\pm$ 474.89	<b>3514.39 <math>\pm</math> 676.80</b>



Figure 4: Training curves of environments from the OpenAI Gym benchmark. Performance is evaluated for 10 episodes over 5 random seeds.

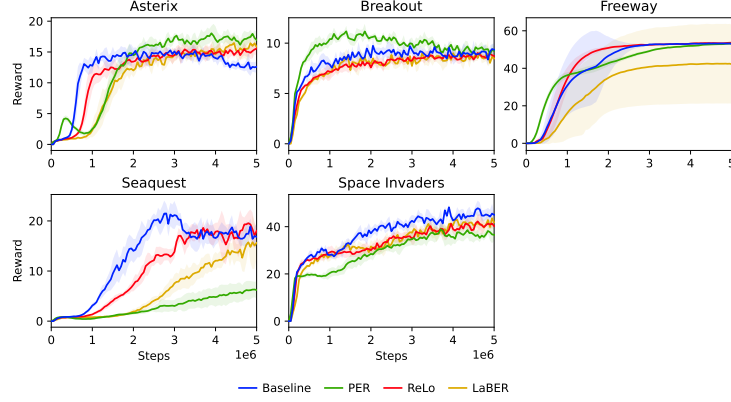


Figure 5: Training curves of environments from the MinAtar benchmark. Performance is evaluated using a running average over the last 1000 episodes over 5 random seeds.

Table 4: Comparison of PER and ReLo on the MinAtar benchmark

	BASLINE	PER	LABER	RELO
ASTERIX	$12.54 \pm 1.08$	<b><math>16.16 \pm 1.02</math></b>	$15.51 \pm 1.11$	$15.68 \pm 0.89$
BREAKOUT	<b><math>9.36 \pm 0.29</math></b>	$8.88 \pm 0.72$	$8.81 \pm 0.61$	$8.98 \pm 0.75$
FREEWAY	$52.80 \pm 0.35$	$52.75 \pm 0.22$	$41.98 \pm 20.99$	<b><math>53.25 \pm 0.37</math></b>
SEAQUEST	$16.13 \pm 2.88$	$6.02 \pm 1.92$	$14.63 \pm 2.98$	<b><math>18.13 \pm 1.25</math></b>
SPACE INVADERS	<b><math>45.36 \pm 1.65</math></b>	$37.36 \pm 4.45$	$43.67 \pm 3.17$	$38.54 \pm 2.60$

## 84 G Arcade Learning Environment

85 We evaluate agents on a compute-constrained version of the Arcade Learning Environment Bellemare  
 86 et al. [2013], training each agent for 2M frames. We chose the standard 24 environments from the  
 87 suite for our evaluation. ReLo is competitive with PER Schaul et al. [2016] in the tested environments.  
 88 The training curves for the Temporal Difference Error and the rewards are given in Fig. 6 & Fig. 7  
 89 respectively.

## 90 H Gridworld

91 We implement a simple GridWorld for the experiments that high-  
 92 lights the drawbacks of PER with TD loss prioritization in Section  
 93 5.6. It consists of a  $7 \times 7$  grid. A visualization of the grid is given  
 94 in Fig. 8. The start state of the agent, represented by the blue point,  
 95 is at the top left and the goal state is at the bottom right, represented  
 96 by the green point. The agent is represented by the black point. The  
 97 agent gets a reward of +2 when it reaches the goal state, but does  
 98 not receive reward anywhere else other than the stochastic point.  
 99 The red point marks the state with a corresponding reward uniformly  
 100 sampled from  $[-0.5, 0.5]$ . The locations of these points is fixed  
 101 and does not change during training or evaluation. The state that  
 102 the agent receives is  $(x, y)$  where  $x$  and  $y$  are the coordinates of  
 103 the agent’s location,  $x, y \in [-3, 3]$  and  $x, y \in \mathbb{Z}$ . The agent has 4  
 104 actions, to move up, down, left or right which will deterministically  
 105 move the agent in that direction by 1 unit. If the agent is at the edge  
 106 of the grid and takes an action that will move it out of the  $7 \times 7$  grid,  
 107 then it remains in the same location.

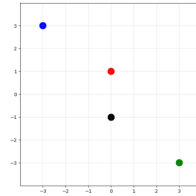


Figure 8: A top down view of the GridWorld. The agent is the black point. It starts at the blue point and the goal state is the green point. The red point represents the location of the stochastic reward.

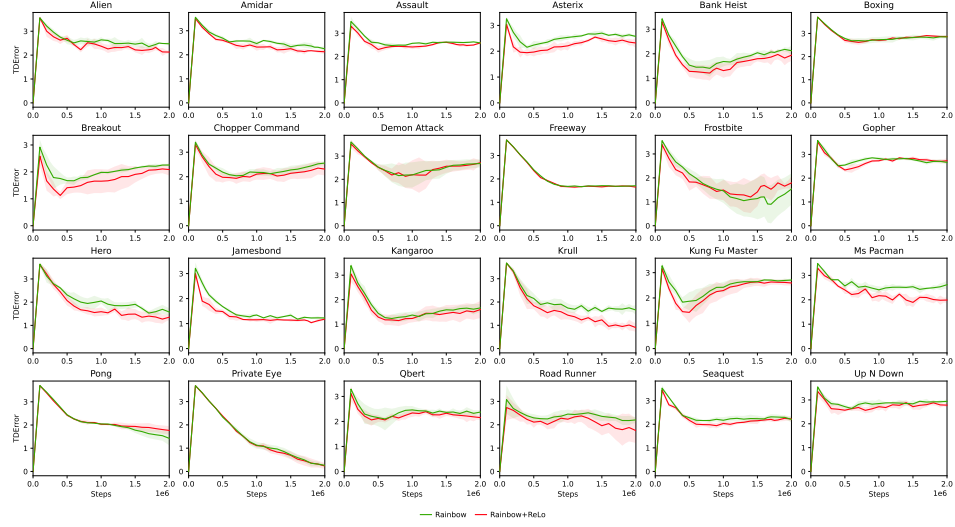


Figure 6: Temporal difference loss curves for Rainbow (with PER) and Rainbow with ReLo. Rainbow with ReLo achieves lower loss compared to PER, showing that ReLo is able to prioritize samples with reducible loss. The dark line represents the mean and the shaded region is the standard deviation over 3 seeds.

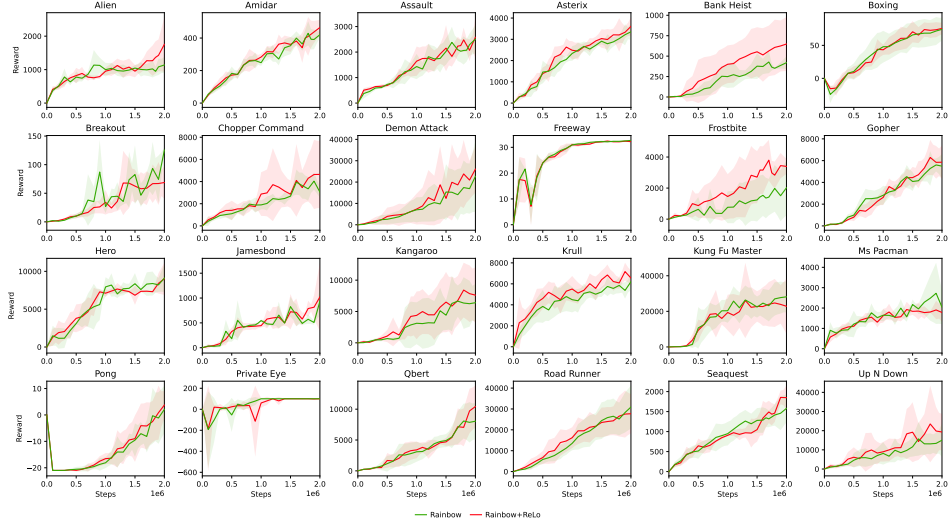


Figure 7: Training curves of 24 environments from the ALE benchmark. Performance is evaluated for 10 episodes over 5 random seeds.

## References

- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47: 253–279, 2013.
- Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning*, pages 449–458. PMLR, 2017.
- Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Matteo Hessel, Ian Osband, Alex Graves, Volodymyr Mnih, Rémi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. Noisy networks for exploration. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=rywHCPkAW>.

Table 5: Comparison of Rainbow with PER and Rainbow with ReLo on the ALE benchmark

	RAINBOW	RAINBOW W/ ReLo
ALIEN	1278.70 $\pm$ 223.14	<b>1352.90 <math>\pm</math> 535.70</b>
AMIDAR	376.20 $\pm$ 81.07	<b>410.10 <math>\pm</math> 85.63</b>
ASSAULT	2241.78 $\pm$ 648.17	<b>2617.37 <math>\pm</math> 555.97</b>
ASTERIX	3214.50 $\pm$ 323.70	<b>3352.00 <math>\pm</math> 431.98</b>
BANKHEIST	526.80 $\pm$ 277.83	<b>641.80 <math>\pm</math> 284.67</b>
BOXING	<b>76.65 <math>\pm</math> 14.58</b>	76.21 $\pm$ 11.59
BREAKOUT	<b>82.03 <math>\pm</math> 46.10</b>	68.36 $\pm$ 45.39
CHOPPERCOMMAND	2794.00 $\pm$ 732.87	<b>4974.00 <math>\pm</math> 2801.93</b>
DEMONATTACK	25500.50 $\pm$ 16311.26	<b>29294.30 <math>\pm</math> 15905.25</b>
FREEWAY	<b>32.58 <math>\pm</math> 0.31</b>	32.35 $\pm$ 0.25
FROSTBITE	2850.60 $\pm$ 1553.34	<b>3532.30 <math>\pm</math> 1270.73</b>
GOPHER	5336.60 $\pm$ 1023.75	<b>5679.80 <math>\pm</math> 1199.68</b>
HERO	<b>9907.15 <math>\pm</math> 2108.71</b>	8830.30 $\pm$ 1894.51
JAMESBOND	810.00 $\pm$ 412.00	<b>902.00 <math>\pm</math> 469.40</b>
KANGAROO	<b>8904.00 <math>\pm</math> 3879.97</b>	8091.00 $\pm$ 3618.47
KRULL	6553.18 $\pm$ 803.15	<b>6718.67 <math>\pm</math> 799.47</b>
KUNGFUMASTER	<b>29371.00 <math>\pm</math> 8525.69</b>	23654.00 $\pm$ 12360.40
MSPACMAN	<b>2094.70 <math>\pm</math> 614.58</b>	1755.80 $\pm$ 374.32
PONG	3.18 $\pm$ 9.02	<b>5.96 <math>\pm</math> 6.45</b>
PRIVATEEYE	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00
QBERT	8382.00 $\pm$ 2935.16	<b>10900.25 <math>\pm</math> 2704.23</b>
ROADRUNNER	<b>29333.00 <math>\pm</math> 9465.78</b>	29222.00 $\pm$ 8696.91
SEAQUEST	1377.20 $\pm$ 362.57	<b>1848.80 <math>\pm</math> 788.85</b>
UPNDOWN	17065.40 $\pm$ 7637.25	<b>21241.50 <math>\pm</math> 8599.97</b>

- 119 Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error  
120 in actor-critic methods. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th*  
121 *International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning*  
122 *Research*, pages 1587–1596. PMLR, 10-15 Jul 2018. URL <https://proceedings.mlr.press/v80/fujimoto18a.html>.  
123
- 124 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
125 maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer G. Dy and  
126 Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*,  
127 *ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings*  
128 *of Machine Learning Research*, pages 1856–1865. PMLR, 2018. URL <http://proceedings.mlr.press/v80/haarnoja18b.html>.  
129
- 130 Matteo Hessel, Joseph Modayil, Hado van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney,  
131 Daniel Horgan, Bilal Piot, Mohammad Gheshlaghi Azar, and David Silver. Rainbow: Combining  
132 improvements in deep reinforcement learning. *CoRR*, abs/1710.02298, 2017. URL <http://arxiv.org/abs/1710.02298>.  
133
- 134 Zhang-Wei Hong, Tao Chen, Yen-Chen Lin, J. Pajarinen, and Pulkit Agrawal. Topological experience  
135 replay. *ICLR*, 2022. doi: 10.48550/arXiv.2203.15845.
- 136 Yuenan Hou, Lifeng Liu, Qing Wei, Xudong Xu, and Chunlin Chen. A novel ddpg method with  
137 prioritized experience replay. In *Systems, Man, and Cybernetics (SMC), 2017 IEEE International*  
138 *Conference on*, pages 316–321. IEEE, 2017.
- 139 Aviral Kumar, Abhishek Gupta, and Sergey Levine. Discor: Corrective feedback in reinforcement  
140 learning via distribution correction. *Advances in Neural Information Processing Systems*, 33:  
141 18560–18572, 2020.
- 142 Thibault Lahire, M. Geist, and E. Rachelson. Large batch experience replay. *ICML*, 2021.
- 143 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,  
144 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In Yoshua

- 145 Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations,*  
 146 *ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL  
 147 <http://arxiv.org/abs/1509.02971>.
- 148 Xu-Hui Liu, Zhenghai Xue, Jing-Cheng Pang, Shengyi Jiang, Feng Xu, and Yang Yu. Regret  
 149 minimization experience replay in off-policy reinforcement learning. *NEURIPS*, 2021.
- 150 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G.  
 151 Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Pe-  
 152 tersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumaran,  
 153 Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep rein-  
 154 forcement learning. *Nature*, 518(7540):529–533, February 2015. ISSN 00280836. URL  
 155 <http://dx.doi.org/10.1038/nature14236>.
- 156 Youngmin Oh, Jinwoo Shin, Eunho Yang, and Sung Ju Hwang. Model-augmented prioritized  
 157 experience replay. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=WuEiafqdy9H>.
- 158 Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. In  
 159 Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representa-*  
 160 *tions, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.  
 161 URL <http://arxiv.org/abs/1511.05952>.
- 162 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
 163 optimization algorithms. *arXiv preprint arXiv: Arxiv-1707.06347*, 2017.
- 164 Samarth Sinha, Jiaming Song, Animesh Garg, and S. Ermon. Experience replay with likelihood-free  
 165 importance weights. *LADC*, 2020.
- 166 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden,  
 167 Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Riedmiller.  
 168 Deepmind control suite. *arXiv preprint arXiv: Arxiv-1801.00690*, 2018.
- 169 Hado van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-  
 170 learning. In Dale Schuurmans and Michael P. Wellman, editors, *Proceedings of the Thirtieth*  
 171 *AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*, pages  
 172 2094–2100. AAAI Press, 2016. URL [http://www.aaai.org/ocs/index.php/AAAI/AAAI16/](http://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/12389)  
 173 [paper/view/12389](http://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/12389).
- 174 Che Wang and Keith Ross. Boosting soft actor-critic: Emphasizing recent experience without  
 175 forgetting the past. *arXiv preprint arXiv: Arxiv-1906.04009*, 2019.
- 176 Ziyu Wang, Tom Schaul, Matteo Hessel, Hado van Hasselt, Marc Lanctot, and Nando de Freitas.  
 177 Dueling network architectures for deep reinforcement learning. In Maria-Florina Balcan and  
 178 Kilian Q. Weinberger, editors, *Proceedings of the 33rd International Conference on Machine*  
 179 *Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop*  
 180 *and Conference Proceedings*, pages 1995–2003. JMLR.org, 2016. URL [http://proceedings.](http://proceedings.mlr.press/v48/wangf16.html)  
 181 [mlr.press/v48/wangf16.html](http://proceedings.mlr.press/v48/wangf16.html).
- 182 Denis Yarats and Ilya Kostrikov. Soft actor-critic (sac) implementation in pytorch. [https://github.](https://github.com/denisyarats/pytorch_sac)  
 183 [com/denisyarats/pytorch\\_sac](https://github.com/denisyarats/pytorch_sac), 2020.
- 184 Kenny Young and Tian Tian. Minatar: An atari-inspired testbed for thorough and reproducible  
 185 reinforcement learning experiments. *arXiv preprint arXiv:1903.03176*, 2019.
- 186 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse  
 187 reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.