# Just Round: Quantized Observation Spaces Enable Memory Efficient Learning of Dynamic Locomotion

Lev Grossman[1] and Brian Plancher[2]

*Abstract*— **Deep reinforcement learning (DRL) is one of the most powerful tools for synthesizing complex robotic behaviors. But training DRL models is incredibly compute and memory intensive, requiring large training datasets and replay buffers to achieve performant results. This poses a challenge for the next generation of field robots that will need to learn on the edge to adapt to their environment. In this paper, we begin to address this issue through observation space quantization. We evaluate our approach using four simulated robot locomotion tasks and two state-of-the-art DRL algorithms, the on-policy Proximal Policy Optimization (PPO) and off-policy Soft Actor-Critic (SAC) and find that observation space quantization reduces overall memory costs by as much as $4.2\times$ without impacting learning performance.**

## I. INTRODUCTION

Deep reinforcement learning (DRL) continues to see increased attention by the robotics community due to its ability to learn complex behaviors in both simulated and real environments. These methods have been successfully applied to a host of robotic tasks including: dexterous manipulation [1], quadrupedal locomotion [2], and high-speed drone racing [3]. Despite these successes, DRL remains largely sample inefficient, depending on enormous amounts of training data to learn. As much of this data is kept in replay buffers during training, DRL is extremely memory intensive, limiting the number of computational platforms that can support such operations, and largely confining state-of-the-art model training to the cloud. For example, OpenAI used 400 compute devices and collected two years worth of experience data per hour in simulation in order to train a physical dexterous robot hand to solve a Rubik's cube [4].

At the same time, there has been a recent push to enable learning on physical robot hardware [5]. Using environment experience collected directly on-device and sending this batched data over Ethernet to a workstation, researchers have shown that DRL, and off-policy methods in particular, hold promise in learning robust locomotive policies on physical robot hardware [6], [7].

Unfortunately, for such real-world robots, a virtually unlimited compute and memory budget is unrealistic. Still, many robots, especially those involved in tasks as consequential as search-and-rescue and space exploration [8], will have to adapt to ever-changing environmental conditions and continue to optimize and update their internal policies over the course of their lifetime [9]. As such, to untether these

[1]Lev Grossman is with Berkshire Grey, Bedford, MA, USA. lev.grossman@berkshiregrey.com
[2]Brian Plancher is with Barnard College, Columbia University, New York, NY, USA. bplancher@barnard.edu

methods from the confines of a lab, data collection and storage must be memory efficient enough to allow for either low-latency networking [10], [11] or for sufficient experience data to be stored on-board edge computing devices [12]. As fast and secure updates may not be possible in remote locations or when using bandwidth-constrained or high-latency cloud networks [13], it is imperative to find ways to reduce the overall memory footprint of DRL training.

In this work we begin to address this issue through observation space quantization. We focus on the observation space in particular, as the observations stored in a replay buffer generally consume over 90% of the total memory footprint of DRL training for state-of-the-art locomotion policies [2], [14]. Importantly, unlike simply reducing the number of observations stored in the buffer, which decreases the memory footprint at the cost of reduced learning performance, our quantization scheme is able to reduce memory usage without impacting the training performance. We present experiments across four popular simulated robotic locomotion domains, using two of the most popular DRL algorithms, the on-policy Proximal Policy Optimization (PPO) and off-policy Soft Actor-Critic (SAC), and find that our approach can reduce the memory footprint by as much as $4.2\times$ without impacting training performance. We open-source our implementation for use by the wider robot learning community.

## II. RELATED WORK

**Locomotive Learning and Efficiency:** Deep reinforcement learning (DRL) methods have been successfully applied to a variety of complex simulated [15], [16] and real [17] robotic locomotion tasks. Advances in asynchronous algorithms [18] and massively parallel simulation [19] have enabled faster learning. However, DRL algorithms often remain data intensive and prone to being sample inefficient. This inefficiency can be especially impactful when transferring learned policies from simulation to physical robots [20] or when learning policies directly on hardware [6], [7]. To address this, some have found success in learning using substantially reduced buffer sizes and intelligent replay sampling [21]. Others have been able to increase the sample efficiency of DRL based locomotion by using more powerful off-policy algorithms [14]. Still, improving the sample efficiency and thus memory efficiency of DRL remains a major area of interest in the robotic learning community.

**Compressed and embedded spaces:** Previous work on learning in compressed or embedded spaces has mainly surrounded model-based techniques [22], [23], often learning

an embedding directly from images [24], [25]. While model-based learning has been shown to work well in some complex dynamic environments [26], model-free methods remain a popular choice in the dynamic locomotion community [16], [17]. Others have turned to embedded and more descriptive action spaces [27], [28], [29], [30] and reduced order models [23] to enable more robust and sample efficient learning. However, these efforts have mainly ignored the impact of observation space compression on model-free learning.

**Quantization:** Quantization or discretization of deep neural network weights and parameters has seen increased popularity [31], [32]. These methods are able to achieve competitive performance on classic deep computer vision tasks while reducing the size footprint of the convolutional neural network (CNN) models used. More recently, work has been done to adapt these same techniques to reinforcement learning (RL) [33], [34]. This work has mainly focused on the quantization of the parameters of the critic or value networks, in order to reduce overall model size and speed up learning. However, little has been done to evaluate the effect of applying quantization to a task's observation space.

## III. LEARNING BACKGROUND

Reinforcement learning poses problems as Markov decision processes (MDPs), where an MDP is defined by a set of observed and hidden states, $\mathcal{S}$, actions, $\mathcal{A}$, stochastic dynamics, $p(s_{t+1}|s_t, a_t)$, a reward function $r(s, a)$, and a discount factor, $\gamma$. The RL objective is to compute the policy, $\pi^*(s, a)$, that maximizes the expected discounted sum of future rewards, $\mathbb{E}_{s,a}\left(\sum_t \gamma^t r_t\right)$. We train all learning tasks with two state-of-the-art on- and off-policy reinforcement learning algorithms: Proximal Policy Optimization (PPO) and Soft Actor-Critic (SAC).

### A. Proximal Policy Optimization (PPO)

Proximal Policy Optimization [35] is an on-policy, policy gradient [36] algorithm that employs an actor-critic framework to learn both the optimal policy, $\pi^*(s, a)$, as well as the optimal value function, $V^*(s)$. Both the policy $\pi_\theta$ and value function $V_\phi$ are parameterized by neural networks with weights $\theta$ and $\phi$ respectively. Similar to Trust Region Policy Optimization (TRPO) [37], PPO stabilizes policy training by penalizing large policy updates. Additionally, as all updates are computed with samples taken from the current policy, PPO often requires high sample complexity.

During training PPO needs to store the parameters of both its value and policy networks[1]–usually shallow multilayer perceptrons (MLPs)–as well as its on-policy rollout buffer $\mathcal{D}_k$. $\mathcal{D}_k$ stores $(s, a, r)$ tuples, which are refreshed during each iteration of the algorithm.

By combining both the models and rollout buffer, the general space complexity of PPO can be written as:

$$\texttt{size} = \texttt{sizeof}(\pi_\theta) + \texttt{sizeof}(V_\phi) + \texttt{sizeof}(\mathcal{D}_k). \quad (1)$$

Importantly, despite throwing out and re-collecting an entirely new rollout buffer during each iteration, $\mathcal{D}_k$ dominates the overall memory footprint of PPO.

For example, Miki et al. [2] uses a student-teacher approach to bridge the sim-to-real gap and enable robust quadrupedal locomotion in rough, wilderness terrain. This process starts by using PPO to train the teacher model–a 3-layer MLP fed by two smaller encoder networks. Assuming all parameters are stored as 64-bit floats, the three models have a combined size of $\sim$1.3MB. In comparison, the 391-dimensional observation space, 16-dimensional action space, and batch size of 8,300 leads to a rollout buffer of over 27MB of space (assuming 64-bit floats). This means that $\mathcal{D}_k$ accounts for 95% of the total memory footprint. Furthermore, this size is dominated by the stored observations, which account for 96% of the size of $\mathcal{D}_k$ and thus 92% of the total memory footprint.

### B. Soft Actor-Critic (SAC)

Soft Actor-Critic [14], [6] is an off-policy RL algorithm that generally extends soft Q-learning (SQL) [38] and optimizes a "maximum entropy" objective, which promotes exploration according to a temperature parameter, $\alpha$:

$$\mathbb{E}_{(s_t, a_t) \sim \pi}\left[\sum_t \gamma^t r(s_t, a_t) + \alpha \mathcal{H}\left(\pi(\cdot|s_t)\right)\right]. \quad (2)$$

SAC makes a number of improvements on SQL by automatically tuning the temperature parameter $\alpha$ using double Q-learning, similar to the Twin Delayed DDPG (TD3) algorithm [39], to correct for overestimation in the Q-function, and learning not only the Q-functions and the policy, but also the value function.

Like PPO, SAC's memory usage comes from its models and off-policy replay buffer. Shallow MLPs are once again used to approximate the policy $\pi_\theta$ as well as the two Q-functions $Q_{\phi_1}, Q_{\phi_2}$. Unlike PPO, however, the replay buffer, $\mathcal{D}$, is generally much larger in size as it stores trajectories from every iteration, usually acting as a size limited queue.

By combining the models and replay buffer, the general space complexity of SAC can be written as:

$$\texttt{size} = \texttt{sizeof}(\pi_\theta) + 2 * \texttt{sizeof}(Q_\phi) + \texttt{sizeof}(\mathcal{D}). \quad (3)$$

As in the case of PPO, $\mathcal{D}$ dominates the memory footprint of SAC. For example, Haarnoja et al. [14] used SAC to train the quadruped Minitaur to walk. The combined number of parameters in their policy and two value networks (ignoring bias terms and again assuming 64-bit floats) was about 2.3MB. With an 112-dimensional observation space, an 8-dimensional action space, and a replay buffer of size 1e6,[2] $\mathcal{D}$ consumed about 96.8MB of memory, equating to 97.7% of the total memory footprint. Again, the size of $\mathcal{D}$ is dominated by the observations, which account for 92.5% of its size and 90.4% of the total memory footprint.

---

[1]Many PPO implementations save space by using one central MLP with two additional single-layer policy and value function model heads. We focus on the standard PPO model approach in this section for clarity.

[2]A conservative estimate as they collect 100k-200k samples.

## IV. METHOD

In this work, we focus on one particular type of quantization: numerical rounding. We define the function $\mathrm{round}(x, m)$ to be the result of rounding $x \in \mathbb{R}$ to $m \in \mathbb{N}$ decimal places and extend this operation to all real vectors $X \in \mathbb{R}^n$. Thus, rounding an observed state, $\mathrm{round}(s, m)$, to $m$ decimal places is the same as $\{\mathrm{round}(s_i, m), \forall s_i \in s\}$.

We also make use of a `clamp` operation with bounds $a, b$:

$$\mathrm{clamp}(s, a, b) = \min(\max(s, a), b), \qquad (4)$$

This operation is useful for further reducing the number of bits required to express each individual observation $s$ and is a popular tool employed in previous quantization efforts [40].

Combining both operations, we transform an input observation state $s$ into a quantized observation state $s_q$:

$$s_q = \mathrm{round}(\mathrm{clamp}(s, a, b), m). \qquad (5)$$

In practice, we set $a = -b$ in order to clamp symmetrically around 0. We can thus calculate the number of bits required to encode a quantized output as:

$$1 + \lceil \log_2(b) \rceil + \lceil m \log_2(10) \rceil, \qquad (6)$$

where 1 bit is needed for the sign, $\lceil \log_2(b) \rceil$ bits are needed to encode the integer left of the decimal, and $\lceil m \log_2(10) \rceil$ bits are required to encode $m$ decimal places.

We adapted out quantization scheme to the data from our baseline experiments. In those experiments we found that $\max(|s|) < 128$, which allowed us to set our clamping bounds at $a = -127, b = 127$ and use 7 integer bits. Combined with the 1 sign bit and either 4 ($m = 1$) or 7 ($m = 2$) decimal bits, we were able to quantize the original 64-bit floating-point values down to either 12 or 15 bits, an over $5\times$ and $4\times$ reduction in space respectively.

## V. EXPERIMENTS AND RESULTS

We evaluate the effectiveness of our quantization framework across four popular, simulated robotic locomotion environments made available by OpenAI Gym [41]: `Walker2d-v3`, `HalfCheetah-v3`, `Ant-v3`, and `Humanoid-v3`. Figure 1 depicts the four environments using the MuJoCo physics engine [42].

### A. Implementation Details

Our full implementation is publicly available and can be found in our GitHub repository: `https://github.com/A2R-Lab/Just-Round`.

We train all agents using the `stable-baselines3` [43] implementations of both PPO and SAC, modifying only the buffer logic in order to accommodate quantization of the stored observations. Unless otherwise noted, we run all experiments using the default PPO and SAC hyperparameters, only adjusting the number of total training steps.[3]

---

[3]Training done on an Intel Core i7-8700K CPU and an NVIDIA GeForce GTX 1080 Ti GPU.
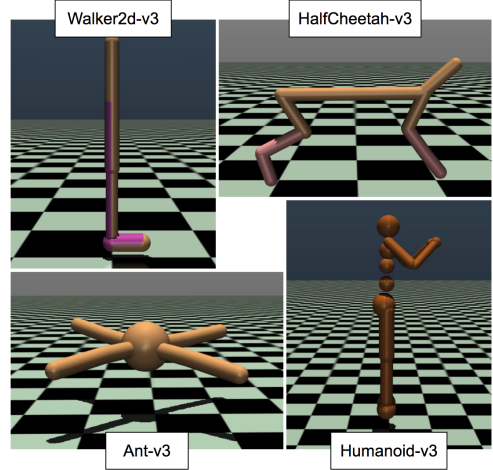


Fig. 1: MuJoCo simulated robotic locomotion environments: `Walker2d-v3` (top-left), `HalfCheetah-v3` (top-right), `Ant-v3` (bottom-left), `Humanoid-v3` (bottom-right).

### B. Memory Efficiency

Figure 2 compares, for all locomotion environments and for both PPO and SAC, the memory required both with and without quantizing the observations in the buffer. We breakdown the memory usage into the the policy and value/Q-value networks (shown in black) and the buffer (shown in blue). As mentioned in Section III, the memory footprint is dominated by the buffers, rendering the model's memory footprint barely visible in the PPO plots and mostly invisible in the SAC plots, where it only accounts for less than 2% of the total memory.

For all experiments we set $a = 127, b = -127$. For PPO we set $m = 1$, and for SAC we set $m = 2$ after observing that SAC required extra decimal precision to achieve optimal convergence. We hypothesize that this discrepancy between PPO and SAC is a result of the off-policy nature of SAC. While in each iteration PPO is trained using a smaller number of rollouts from the most recent policy, SAC utilizes a larger set of replay data collected from both past and present policies. This larger buffer may in turn lead to a higher chance of observation "collisions," where two rounded tuples $(s_q, a, r_1)$ and $(s_q, a, r_2)$ present the same quantized observation and action with a different reward, leading to counterproductive training steps.

Overall, we found that quantizing only the observations stored in the buffers reduces overall memory footprint by $2.1\times$ to $4.2\times$ for PPO and $2.6\times$ to $3.9\times$ for SAC. In the `Humanoid-v3` environment, this $3.9\times$ reduction in memory used for SAC decreases a 3GB memory requirement to just under 750MB. This space savings is significant, allowing real, memory-constrained learning systems the ability to store much larger replay buffers on-board than would be possible without quantization.

### C. Training Performance

Figure 3 plots the learning curves for the first 500k steps of training for the `Humanoid-v3` environment using both quantization and the baseline method of reducing the number
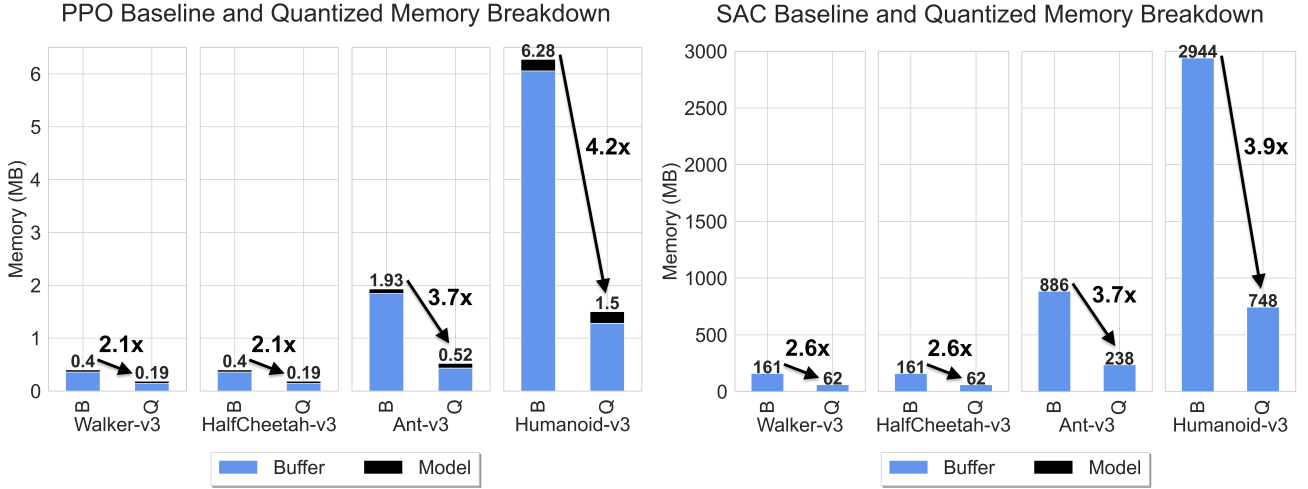
Fig. 2: Memory breakdown of PPO (left) and SAC (right) both without (B) and with (Q) observation space quantization across environments. The model size is not visible in the SAC chart as it comprises less than 2% of the total memory.
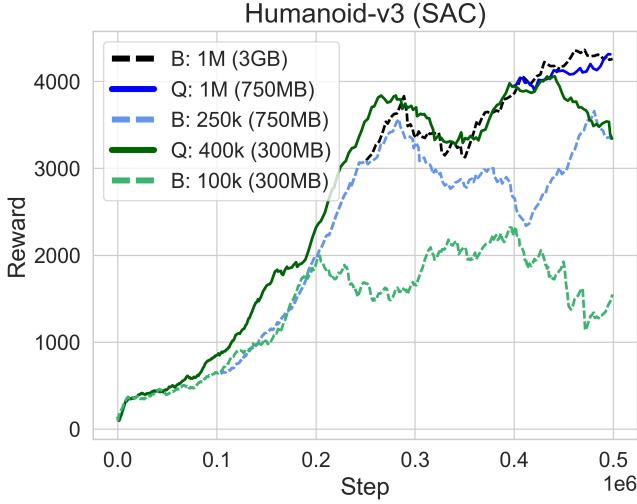


Fig. 3: Learning curves of SAC applied to the `Humanoid-v3` environment with (Q, solid) and without (B, dashed) quantization and with variable replay buffer sizes: 1M, 400k, 250k, 100k. Colors differentiate experiments based on memory footprint: black ~3GB, dark/light blue ~750MB, dark/light green ~300MB. Parameters $m = 2$, $a = 127$, and $b = -127$ found through minimal tuning.

of observations in the replay buffer. We ran all experiments with the same hyperparameters, only varying the buffer quantization and size. The baseline SAC agents with varied numbers of observations in the rollout buffer are plotted using dashed lines, and the quantized agents are plotted with solid lines. We group and color-code experiments by memory footprint: black corresponding to ~3GB in size, dark/light blue to ~750MB, and dark/light green to ~300MB.

First and foremost, we find that the quantized agent with a buffer size of 1M (in solid dark blue) displays similar convergence results as the baseline 1M agent (in dashed black), while only using roughly a quarter of the memory (~750MB vs ~3GB). Holding the memory footprint con-

stant, we find that the quantized agents (solid dark blue and dark green) converge to significantly higher rewards than their unquantized counterparts (dashed light blue and light green). While stark, this difference is not entirely surprising. For more complex learning tasks, like humanoid walking, larger replay buffers are often essential in enabling performant convergence. Overall, this experiment indicates that observation-based quantization may allow us to retain large replay buffers and their associated highly performant training while reducing memory costs by as much as $4.2\times$.

### D. Training Convergence and Speed

To validate the generalizability of this result we show the learning curves of PPO and SAC across four different simulated locomotion environments in Figure 4. These figures report the average of ten random seeds and also display the standard deviation with the accompanying shaded regions.

We find that convergence is not substantially affected by the introduction of observation-based quantization. All final rewards of both quantized and unquantized agents are within a standard deviation of each other, and in some cases the quantized models are able to train a policy with a higher total reward. Furthermore, we show in Table I that the overhead from performing quantization also does not significantly impact overall training times and in some cases actually speeds up overall training through reductions in memory access overheads.

| | | Walker2d | HalfCheetah | Ant | Humanoid |
|---|---|---|---|---|---|
| | B | 1.66 | 1.77 | 2.16 | 1.78 |
| PPO | Q | 1.61 | 1.76 | 2.17 | 1.87 |
| | T | 0.97× | 0.99× | 1.00× | 1.05× |
| | B | 8.38 | 8.49 | 9.00 | 9.83 |
| SAC | Q | 8.41 | 8.44 | 8.97 | 9.91 |
| | T | 1.00× | 0.99× | 1.00× | 1.01× |

TABLE I: Wall clock time in milliseconds per step of PPO and SAC with (Q) and without (B) quantization as well as the relative total time (T) of training with quantization.
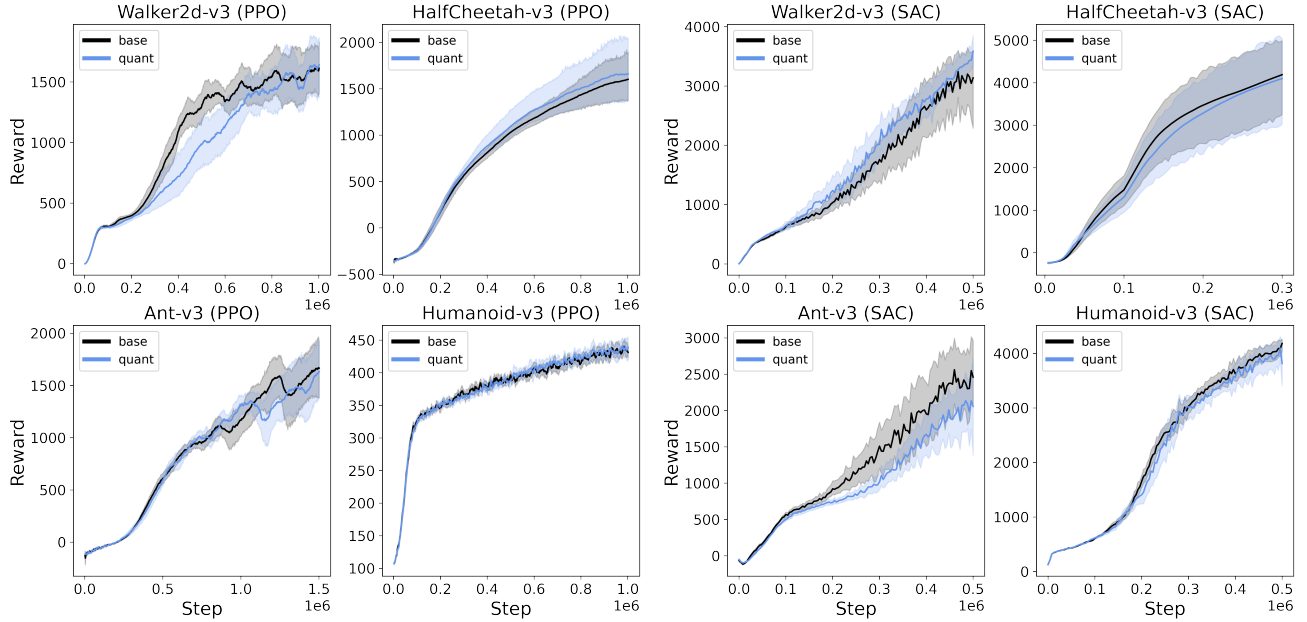
Fig. 4: Learning curves of PPO (left) and SAC (right) with (blue) and without (black) quantization across four simulated locomotion environments. Rewards are averaged across ten trials and standard deviations are shaded.

## VI. CONCLUSION AND FUTURE WORK

This paper presents observation space quantization, a method for reducing the overall memory requirements of DRL without impacting training performance.

We evaluated the overall memory footprint and learning performance for four simulated robotic locomotion tasks using two state-of-the-art on- and off- policy DRL algorithms. We found that applying observation space quantization reduced the effective memory footprint of DRL training by as much as $4.2\times$, while preserving both training convergence and speed. We believe that this approach can be a powerful way to reduce the overall memory consumption of DRL.

While we admit that not every learning-based solution requires memory compression (especially if learning is being done on a powerful workstation with ample storage), if we are to achieve lifelong, practical learning on the edge, curbing memory usage is a must.

In future work, we hope to compare our approach to other forms of reduced precision (e.g., half-precision floating point [34]). We also hope to deploy our approach onto physical robot hardware and test it in the context of real-world edge RL. Finally, we hope to apply our approach in the context of tiny robot learning [44] and help usher in a new era of ubiquitous edge RL.

## APPENDIX A: LEARNING HYPERPARAMETERS

All experiments were carried out using the default `stable-baselines3` learning parameters for PPO and SAC unless otherwise noted. The full hyperparameter list is given in Table II. Both PPO and SAC used two-layered feed-forward networks with 64 nodes per layer.

| | Parameter | Value |
|---|---|---|
| PPO | MDP steps per update | 2048 |
| | batch size | 64 |
| | learning rate | 0.0003 |
| | gamma | 0.99 |
| | $\lambda$ for GAE | 0.95 |
| | value function cost | 0.5 |
| | entropy cost | 0.0 |
| | max gradient norm | 0.5 |
| | cliprange | 0.2 |
| SAC | learning rate | 0.0003 |
| | buffer size | 1000000 |
| | entropy coefficient | auto |
| | MDP steps between updates | 1 |
| | batch size | 256 |
| | tau | 0.005 |
| | gamma | 0.99 |
| | gradient steps per update | 1 |

TABLE II: Learning hyperparameters

## REFERENCES

[1] OpenAI, M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba, "Learning dexterous in-hand manipulation," 2018. [Online]. Available: https://arxiv.org/abs/1808.00177

[2] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning robust perceptive locomotion for quadrupedal robots in the wild," *Science Robotics*, vol. 7, no. 62, jan 2022. [Online]. Available: https://doi.org/10.1126%2Fscirobotics.abk2822

[3] Y. Song, M. Steinweg, E. Kaufmann, and D. Scaramuzza, "Autonomous drone racing with deep reinforcement learning," 2021. [Online]. Available: https://arxiv.org/abs/2103.08624

[4] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. Zhang, "Solving rubik's cube with a robot hand," 2019. [Online]. Available: https://arxiv.org/abs/1910.07113

[5] J. Ibarz, J. Tan, C. Finn, M. Kalakrishnan, P. Pastor, and S. Levine, "How to train your robot with deep reinforcement learning: lessons we have learned," *The International Journal of Robotics Research*, vol. 40, no. 4-5, pp. 698–721, jan 2021. [Online]. Available: https://doi.org/10.1177%2F0278364920987859

[6] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, "Soft Actor-Critic Algorithms and Applications," *arXiv:1812.05905 [cs, stat]*, Dec. 2018.

[7] P. Wu, A. Escontrela, D. Hafner, K. Goldberg, and P. Abbeel, "Daydreamer: World models for physical robot learning," 2022. [Online]. Available: https://arxiv.org/abs/2206.14176

[8] P. R. Gankidi and J. Thangavelautham, "Fpga architecture for deep learning and its application to planetary robotics," in *2017 IEEE Aerospace Conference*, 2017, pp. 1–9.

[9] S. Thrun and T. M. Mitchell, "Lifelong robot learning," *Robotics and Autonomous Systems*, vol. 15, no. 1, pp. 25–46, 1995, the Biology and Technology of Intelligent Autonomous Agents. [Online]. Available: https://www.sciencedirect.com/science/article/pii/092188909500004Y

[10] B. Liu, L. Wang, and M. Liu, "Lifelong federated reinforcement learning: A learning architecture for navigation in cloud robotic systems," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 4555–4562, oct 2019. [Online]. Available: https://doi.org/10.1109%2Flra.2019.2931179

[11] J. Qi, Q. Zhou, L. Lei, and K. Zheng, "Federated reinforcement learning: Techniques, applications, and open challenges," 2021. [Online]. Available: https://arxiv.org/abs/2108.11887

[12] W. Fedus, P. Ramachandran, R. Agarwal, Y. Bengio, H. Larochelle, M. Rowland, and W. Dabney, "Revisiting fundamentals of experience replay," 2020. [Online]. Available: https://arxiv.org/abs/2007.06700

[13] M. Groshev, G. Baldoni, L. Cominardi, A. de la Oliva, and R. Gazda, "Edge robotics: are we ready? an experimental evaluation of current vision and future directions," *Digital Communications and Networks*, 2022. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2352864822000888

[14] T. Haarnoja, S. Ha, A. Zhou, J. Tan, G. Tucker, and S. Levine, "Learning to walk via deep reinforcement learning," 2018. [Online]. Available: https://arxiv.org/abs/1812.11103

[15] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," 2015. [Online]. Available: https://arxiv.org/abs/1509.02971

[16] X. B. Peng, G. Berseth, K. Yin, and M. Van De Panne, "Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 1–13, 2017.

[17] M. Bogdanovic, M. Khadiv, and L. Righetti, "Model-free reinforcement learning for robust locomotion using demonstrations from trajectory optimization," 2021. [Online]. Available: https://arxiv.org/abs/2107.06629

[18] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," 2016. [Online]. Available: https://arxiv.org/abs/1602.01783

[19] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, "Learning to walk in minutes using massively parallel deep reinforcement learning," 2021. [Online]. Available: https://arxiv.org/abs/2109.11978

[20] J. Tan, T. Zhang, E. Coumans, A. Iscen, Y. Bai, D. Hafner, S. Bohez, and V. Vanhoucke, "Sim-to-real: Learning agile locomotion for quadruped robots," 2018. [Online]. Available: https://arxiv.org/abs/1804.10332

[21] Q. Lan, Y. Pan, J. Luo, and A. R. Mahmood, "Memory-efficient reinforcement learning with knowledge consolidation," 2022. [Online]. Available: https://arxiv.org/abs/2205.10868

[22] O. Rybkin, C. Zhu, A. Nagabandi, K. Daniilidis, I. Mordatch, and S. Levine, "Model-based reinforcement learning via latent-space collocation," 2021. [Online]. Available: https://arxiv.org/abs/2106.13229

[23] R. T. Fawcett, K. Afsari, A. D. Ames, and K. A. Hamed, "Toward a data-driven template model for quadrupedal locomotion," *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 7636–7643, 2022.

[24] M. Watter, J. T. Springenberg, J. Boedecker, and M. Riedmiller, "Embed to control: A locally linear latent dynamics model for control from raw images," 2015. [Online]. Available: https://arxiv.org/abs/1506.07365

[25] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson, "Learning latent dynamics for planning from pixels," 2018. [Online]. Available: https://arxiv.org/abs/1811.04551

[26] A. Nagabandi, G. Kahn, R. S. Fearing, and S. Levine, "Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 7559–7566.

[27] P. Varin, L. Grossman, and S. Kuindersma, "A comparison of action spaces for learning manipulation tasks," 2019. [Online]. Available: https://arxiv.org/abs/1908.08659

[28] A. Allshire, R. Martín-Martín, C. Lin, S. Manuel, S. Savarese, and A. Garg, "Laser: Learning a latent action space for efficient reinforcement learning," 2021. [Online]. Available: https://arxiv.org/abs/2103.15793

[29] S. Karamcheti, A. J. Zhai, D. P. Losey, and D. Sadigh, "Learning visually guided latent actions for assistive teleoperation," 2021. [Online]. Available: https://arxiv.org/abs/2105.00580

[30] J. Wong, V. Makoviychuk, A. Anandkumar, and Y. Zhu, "Oscar: Data-driven operational space control for adaptive and robust robot manipulation," 2021. [Online]. Available: https://arxiv.org/abs/2110.00704

[31] Y. Gong, L. Liu, M. Yang, and L. Bourdev, "Compressing deep convolutional networks using vector quantization," 2014. [Online]. Available: https://arxiv.org/abs/1412.6115

[32] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Quantized neural networks: Training neural networks with low precision weights and activations," 2016. [Online]. Available: https://arxiv.org/abs/1609.07061

[33] S. Krishnan, M. Lam, S. Chitlangia, Z. Wan, G. Barth-Maron, A. Faust, and V. J. Reddi, "Quarl: Quantization for sustainable reinforcement learning," 2019. [Online]. Available: https://arxiv.org/abs/1910.01055

[34] J. Bjorck, X. Chen, C. De Sa, C. P. Gomes, and K. Q. Weinberger, "Low-precision reinforcement learning: Running soft actor-critic in half precision," 2021. [Online]. Available: https://arxiv.org/abs/2102.13565

[35] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," *arXiv:1707.06347 [cs]*, July 2017.

[36] R. S. Sutton, D. M. A. llester, S. Singh, and Y. Mansour, "Policy Gradient Methods for Reinforcement Learning with Function Approximation," 2000, pp. 1057–1063.

[37] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel, "Trust Region Policy Optimization," *arXiv:1502.05477 [cs]*, Feb. 2015.

[38] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, "Reinforcement Learning with Deep Energy-Based Policies," *arXiv:1702.08165 [cs]*, Feb. 2017.

[39] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing Function Approximation Error in Actor-Critic Methods," *arXiv:1802.09477 [cs, stat]*, Feb. 2018.

[40] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. Howard, H. Adam, and D. Kalenichenko, "Quantization and training of neural networks for efficient integer-arithmetic-only inference," 2017. [Online]. Available: https://arxiv.org/abs/1712.05877

[41] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," 2016. [Online]. Available: https://arxiv.org/abs/1606.01540

[42] E. Todorov, T. Erez, Y. Tassa, and tassa, "MuJoCo: A physics engine for model-based control," in *Proceedings of the IEEERAS International Conference on Intelligent Robots*, 2012.

[43] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann, "Stable-baselines3: Reliable reinforcement learning implementations," *Journal of Machine Learning Research*, vol. 22, no. 268, pp. 1–8, 2021. [Online]. Available: http://jmlr.org/papers/v22/20-1364.html

[44] S. M. Neuman, B. Plancher, B. P. Duisterhof, S. Krishnan, C. Banbury, M. Mazumder, S. Prakash, J. Jabbour, A. Faust, G. C. de Croon, and V. Janapa Reddi, "Tiny robot learning: Challenges and directions for machine learning in resource-constrained robots," in *2022 IEEE International Conference on Artificial Intelligence Circuits and Systems (AICAS)*, Incheon, Korea, June 2022.