# The Impact of Quantization and Pruning on Deep Reinforcement Learning Models

Heng Lu[1], Mehdi Alemi[2,3], and Reza Rawassizadeh[1]

[1]Department of Computer Science at Metropolitan College, Boston University, Boston, MA, USA
[2]Department of Orthopaedic Surgery, Harvard Medical School, Boston, MA, USA.
[3]Training Services, MathWorks, Natick, MA, USA.

**Abstract**

Deep reinforcement learning (DRL) has achieved remarkable success across various domains, such as video games, robotics, and, recently, large language models. However, the computational costs and memory requirements of DRL models often limit their deployment in resource-constrained environments. The challenge underscores the urgent need to explore neural network compression methods to make RDL models more practical and broadly applicable. Our study investigates the impact of two prominent compression methods, *quantization* and *pruning* on DRL models. We examine how these techniques influence four performance factors: average return, memory, inference time, and battery utilization across various DRL algorithms and environments. Despite the decrease in model size, we identify that these compression techniques generally do not improve the energy efficiency of DRL models, but the model size decreases. We provide insights into the trade-offs between model compression and DRL performance, offering guidelines for deploying efficient DRL models in resource-constrained settings.

## 1 Introduction and Background

Reinforcement learning has been applied in many fields, including robotics, video games, and recently Reinforcement Learning with Human Feedback (RLHF) [1, 2, 3, 4, 5], has become common in large language models. RLHF methods mitigate biases inherent in language models themselves [4, 6, 7]. Reinforcement learning models that address real-world problems predominantly utilize continuous models based on neural network architecture, known as Deep Reinforcement Learning (DRL).

DRL methods typically involve a world model, agents interacting with the world, and a reward function that evaluates the effectiveness of actions based on the agent's policy towards predefined objectives [8]. Depending on whether the algorithm learns a specific world model, DRL algorithms are categorized into model-based DRL algorithms and model-free DRL algorithms [9]. Model-free DRL methods generally fall into three main categories: deep Q-learning methods [10, 11, 12], policy gradient methods [13, 14, 15], and actor-critic methods [16, 17, 18, 19]. Unlike model-based algorithms, model-free approaches circumvent model bias and offer greater generalizability, which contributes to their popularity in RLHF applications [4, 20].

Neural networks, which are the backbone of DRL methods, are associated with high computational costs and, therefore, resource intensive. Recently, there has been a significant increase in the energy and water consumption of artificial intelligence (AI) data centers [1][2][3][4]. This trend has led to several research studies [21, 22, 23] investigating the resource costs of recent advances in AI. Reducing the energy utilization of DRL has become a crucial need.

Additionally, many systems that benefit from reinforcement learning operate on battery-powered devices, such as extended reality devices and mobile robots. As the size of these devices decreases, their computational capabilities also diminish [24]. One common approach to reducing the computational costs of neural network models is compressing them via pruning and quantization[25, 26]. Network compression methods have been widely applied in computer vision [27, 28] and large language models [29, 30, 31] to improve inference times and reduce memory requirement with minimal compromise to accuracy. Compression enables advanced DRL models to be deployed in robots with low latency and high energy efficiency under constrained resources. Despite its promise, neural network compression in DRL models has received less research attention [32, 33] compared to other fields, such as computer vision and natural language processing.

---

[1]https://www.theatlantic.com/technology/archive/2024/03/ai-water-climate-microsoft/677602
[2]https://www.oregonlive.com/silicon-forest/2022/12/googles-water-use-is-soaring-in-the-dalles-records-show-with-two-more-data-centers-to-come.html
[3]https://www.bloomberg.com/news/articles/2023-07-26/thames-water-considers-restricting-flow-to-london-data-centers
[4]https://www.washingtonpost.com/business/2024/03/07/ai-data-centers-power

Several general approaches have been proposed for neural network compression method quantization and pruning [34, 35]. As for quantization, DRL researchers might be more familiar with vector quantization, which aims to discretize continuous space into a discrete vector set to reduce dimensions [36, 37, 38, 39]. However, in this work, we specifically apply neural network quantization, which focuses on converting float32 format weights and biases into smaller-scale numbers such as int8 or 4-bits rather than vector quantization. Pruning methods, on the other hand, involve removing neurons deemed least important based on criteria such as weight or activation value [40, 41, 42, 43, 44].

In addition to conventional compression approaches [35, 25, 26], there are promising DRL-specific compression approaches [32, 36, 40, 43]. AQuaDem [39] discretizes the action space and learns a discrete set of actions for each state $s$ using behavior cloning from expert demonstrations [45]. The adaptive state aggregation algorithm[36] adaptively discretizes the state space based on Bregman divergence, enabling distinct partitions of the state space. Another group of methods focuses on scalar quantization that reduces the numeric precision of values [46, 40, 33, 32]. NNC-DRL[46] accelerates the DRL training process by speeding up prediction based on GA3C[47] and employs policy distillation to compress the behavior policy network prediction with minimal performance degradation. FIXAR [33] proposes quantization-aware training in fixed points to reduce model size without significant accuracy loss. ActorQ [32] introduces an int8 quantized policy block for rollouts within a traditional distributed RL training loop. PoPS[41] accelerates model speed by initially training a sparse network from a large-scale teacher network through iterative policy pruning, then compacting it into a dense network with minimal performance loss. UVNQ[40] integrates sparse variational dropout[48] with quantization, adjusting the dropout rate to enhance quantization awareness. Dynamic Structured Pruning[43] enhances DRL training by applying a neuron-importance group sparse regularizer and dynamically pruning insignificant neurons based on a threshold. Double Sparse Deep Reinforcement Learning [44] uses multilayer sparse-coding structural network with a nonconvex log regularizer to enforce sparsity while maintaining performance.

In this work, we apply common neural network compression methods, including common quantization and pruning, to five popular deep reinforcement learning models (TRPO[14], PPO[15], DDPG[17], TD3[18], and SAC[19]). We then measure the performance of these algorithms post-compression using metrics such as average return, inference time, and energy usage. In particular, we apply $L_1$ and $L_2$ pruning techniques to these models. For quantization, we utilize *int8* quantization and apply (i) post-training dynamic quantization, (ii) post-training static quantization, and (iii) quantization aware training across the listed models.

To our knowledge, this study represents the first comprehensive evaluation of the effects of pruning and quantization across a range of deep reinforcement learning models. Our experiments and findings offer valuable insights to researchers and developers, assisting them in making informed decisions when choosing between quantization or pruning methods for DRL models. Another prevalent approach for compressing neural network is knowledge distillation [49], but due to its model or application-specific nature (e.g., image classification), we did not include knowledge distillation in our experimental setup.

# 2 Methods

We have applied two types of neural network compression techniques —pruning and quantization— across five prominent DRL models: TRPO[14], PPO[15], DDPG[17], TD3[18], and SAC[19]. This section outlines our quantization methods, followed by our pruning methods.

## 2.1 Quantization

We applied linear quantization across all models, where the relationship between the original input $r$ and its quantized version $q$ is defined as $r = S(q + Z)$. Here, $Z$ represents the zero point in the quantization space, and the scaling factor $S$ maps floating-point numbers to the quantization space. For Post-Training Dynamic Quantization (PTDQ) and Post-Training Static Quantization (PTSQ), we computed $S$ and $Z$ for activations exclusively. In PTSQ, first, baseline models go through a calibration process to compute these quantization parameters and then the models make inferences based on the fixed quantization parameters. In PTDQ, the quantization parameters are computed dynamically. In Quantization-Aware Training (QAT), baseline models are pseudo-quantized during training, meaning computations are conducted in floating-point precision but rounded to integer values to simulate quantization. Subsequently, the original models are converted into quantized versions, and the quantization parameters are stabilized.

## 2.2 Pruning

Neural network pruning typically involves removing neurons within layers, and dependencies can exist where pruning in one layer affects subsequent related layers. The DepGraph approach we employed [50], addresses these dependencies by grouping layers based on their inter-dependencies rather than manually resolving dependencies.

Conceptually, one might consider constructing a grouping matrix $G \in R^{L \times L}$, where $G_{ij} = 1$ signifies a dependency between layer $i$ and layer $j$. However, due to the complexity arising from non-local relations, $G$ can not be easily constructed. Thus, dependency graph $D$ is proposed, which only contains the local dependency between adjacent layers and from which the grouping matrix can be reduced. These dependencies are categorized into two types: inter-layer dependencies, where the output of one layer $i$ connects to the input of another layer $j$, and intra-layer dependencies, such as within BatchNorm layers, where inputs and outputs share the same pruning scheme.

After constructing the dependency graph and determining grouping parameters based on this graph, we utilized a norm-based importance score. However, directly summing importance scores across different layers can lead to meaningless results and potential divergence. Therefore, for a parameter $w$ in group $g$ with $K$ prune-able dimensions, a regularization term $R(g, k)$ is used in sparse training to select the optimal input variables, $R(g, k) = \sum_{k=1}^{K} \gamma_k \cdot I_{g,k}$, where $I_{g,k} = \sum_{w \in g} ||w[k]||_2^2$ is the importance for dimension $k$ in $L_2$ pruning and $\gamma_k = 2^{\alpha(I_g^{max} - I_{g,k})/(I_g^{max} - I_g^{min})}$.

# 3   Experiments

## 3.1   Experimental Settings

Our experiments are structured into two main components: quantization and pruning of DRL algorithms. We evaluated the performance of TRPO[14], PPO[15], DDPG[17], TD3[18], and SAC[19] across five Gymnasium[51] (formerly OpenAI Gym[52]) Mujoco environments including: HalfCheetah, HumanoidStandup, Ant, Humanoid, and Hopper. These models are trained using Gymnasium (formerly OpenAI Gym) environments.

To ensure consistency among our reported results, each experiment has been repeated at least 10 times in the same configuration.

**Quantization and Pruning Libraries:** The implementations of quantization and pruning in neural network libraries are not as mature as other functionalities. For instance, in pyTorch pruning does not remove neurons but merely masks them. To ensure the reliability of our experiments, we evaluated various quantization and pruning libraries and selected those that offer the highest accuracy and resource efficiency.

Therefore, to implement pruning, we explored PyTorch[5] and Torch-pruning. In our experiment, Torch-pruning[6], integrated with DepGraph [50], performed exceptionally well, and thus, we utilized it for pruning purposes. Regarding quantization, we experimented with Pytroch, TensorFlow, and ONNX Runtime[7]. Ultimately, we chose PyTorch for QAT, and ONNX Runtime for PTDQ and PTSQ.

**Hardware Settings:** Our hardware infrastructure included two NVidia RTX 4090 GPUs with 24GB of VRAM, 256GB of RAM, and an Intel Core i9 CPU running at 3.30 GHz. The operating system was Ubuntu 20.04 LTS, and we used CUDA Version 12.0 for GPU operations.

## 3.2   Quantization

To implement quantization, we experimented with three approaches: PTDQ, PTSQ and QAT. Quantization-aware training (QAT) involved initially training quantized models with an equivalent dataset size as the baseline models, followed by exporting them into ONNX runtime for comparative analysis.

### 3.2.1   Average Return

The impact of quantization on average return is reported in Table 1. The table underscores the variability of quantization outcomes across different environments and DRL models. For instance, QAT demonstrates its highest efficacy in HumanoidStandup environments, resulting in improved average returns across models except for PPO. The SAC algorithm generally benefits more from QAT, except in the Hopper environment, where its effectiveness is limited. Overall, PTDQ exhibits superior performance, while PTSQ consistently shows the lowest results. The observed performance discrepancies may stem from distribution shifts between data used for optimal path calculations and that utilized during the calibration phase, which are challenging to rectify due to the stochastic nature of the environment.

---

[5]https://pyTorch.org
[6]https://github.com/VainF/Torch-Pruning
[7]https://onnxruntime.ai

|  |  | Baseline | PTDQ | PTSQ | QAT |
|---|---|---|---|---|---|
| TRPO | HalfCheetah | 978.48 | **1004.65** | 909.52 | 287.94 |
|  | HumanoidStandup | 39444.51 | 37002.02 | 35779.53 | **52252.48** |
|  | Ant | 851.33 | 799.46 | **1055.39** | 970.59 |
|  | Humanoid | 74.49 | 75.09 | 74.98 | **178.36** |
|  | Hopper | 164.06 | **163.84** | 162.93 | 7.49 |
| PPO | HalfCheetah | 1542.98 | **1508.82** | 1482.4 | 432.21 |
|  | HumanoidStandup | 117138.8 | 126509.58 | **128155.96** | 28270.98 |
|  | Ant | 1335.55 | 1493.93 | **1528.36** | 963.21 |
|  | Humanoid | 453.56 | 430.5 | **495.0** | 295.75 |
|  | Hopper | 8.33 | 9.41 | 19.28 | **92.32** |
| DDPG | HalfCheetah | 4475.62 | **4656.76** | 3801.11 | 937.53 |
|  | HumanoidStandup | 82747.99 | 82747.99 | 87346.75 | **115325.59** |
|  | Ant | 589.46 | 630.37 | **896.44** | 540.7 |
|  | Humanoid | 1544.76 | 314.51 | **433.15** | 411.92 |
|  | Hopper | 1419.69 | **1260.39** | 349.31 | 996.61 |
| TD3 | HalfCheetah | 8333.37 | **5204.59** | 3915.44 | 4169.07 |
|  | HumanoidStandup | 77140.94 | 77140.94 | 78492.56 | **82211.28** |
|  | Ant | 3423.51 | **2789.51** | 2728.76 | 1833.74 |
|  | Humanoid | 5035.79 | **441.34** | 287.35 | 80.67 |
|  | Hopper | 3596.54 | **3532.45** | 2842.55 | 1826.21 |
| SAC | HalfCheetah | 10460.06 | 3104.33 | 1805.41 | **6163.8** |
|  | HumanoidStandup | 151015.38 | 109714.88 | 83898.23 | **151213.82** |
|  | Ant | 4021.96 | 2474.88 | 1144.59 | **3119.85** |
|  | Humanoid | 4287.29 | -84.51 | 23.4 | **311.68** |
|  | Hopper | 2539.05 | **3621.34** | 2525.01 | 2998.79 |

Table 1: Average returns of quantization for TRPO, PPO, DDPG, TD3, and SAC. The best quantized version for each DRL model on the specific environments is shown in bold.

### 3.2.2 Resource Utilization

To assess the impact of quantization on resource utilization, we conducted measurements and comparisons of memory usage, inference time, and energy consumption between baseline models and their quantized counterparts. Figure 1 illustrates the differences observed in inference time and energy usage between baseline and quantized models.



Figure 1: Inference time (in seconds), energy usage (in Joules) and memory utilization (in MegaByte) of quantization models.

## 3.3 Pruning

To implement pruning, we utilized the torch-pruning package [8] for all our experiments. Each baseline model underwent $L_1$ and $L_2$ pruning, with various pruning percentages ranging from 5% to 70%. In particular, experimented pruning percentages are as follows: {5%,10%,15%,20%,25%,30%,35%,40%,45%,50%,55%,60%,65%,70%}.

| | | Baseline | Pruned | $L_1$ Pruning Percentage | $L_2$ Pruning Percentage |
|---|---|---|---|---|---|
| TRPO | HalfCheetah | 1003.37 | 905.94 | 0.05 | 0.05 |
| | HumanoidStandup | 35281.42 | 42314.62 | 0.55 | 0.6 |
| | Ant | 768.97 | 979.42 | 0.7 | 0.7 |
| | Humanoid | 74.49 | 79.81 | 0.25 | 0.35 |
| | Hopper | 164.06 | 163.33 | 0.15 | 0.45 |
| PPO | HalfCheetah | 1529.18 | 1453.9 | 0.1 | 0.05 |
| | HumanoidStandup | 128722.71 | 117239.65 | 0.1 | 0.1 |
| | Ant | 1563.54 | 263.37 | 0.05 | 0.05 |
| | Humanoid | 474.86 | 417.99 | 0.1 | 0.15 |
| | Hopper | 21.06 | 7.6 | 0.3 | 0.3 |
| DDPG | HalfCheetah | 5104.06 | 4026.46 | 0.05 | 0.05 |
| | HumanoidStandup | 82747.99 | 83513.05 | 0.7 | 0.7 |
| | Ant | 935.72 | 761.55 | 0.1 | 0.05 |
| | Humanoid | 1659.59 | 1059.68 | 0.05 | 0.05 |
| | Hopper | 1574.66 | 1423.7 | 0.05 | 0.05 |
| TD3 | HalfCheetah | 8298.95 | 6535.42 | 0.05 | 0.05 |
| | HumanoidStandup | 77140.94 | 97061.37 | 0.7 | 0.7 |
| | Ant | 3381.72 | 2564.85 | 0.05 | 0.05 |
| | Humanoid | 5040.01 | 5046.25 | 0.1 | 0.1 |
| | Hopper | 3593.0 | 3589.47 | 0.05 | 0.05 |
| SAC | HalfCheetah | 10467.97 | 10531.88 | 0.05 | 0.05 |
| | HumanoidStandup | 136900.13 | 137574.71 | 0.7 | 0.7 |
| | Ant | 3499.52 | 282.28 | 0.05 | 0.05 |
| | Humanoid | 4251.2 | 3549.46 | 0.25 | 0.35 |
| | Hopper | 2549.96 | 2412.92 | 0.05 | 0.2 |

Table 2: Average returns of pruning results for TRPO, PPO, DDPG, TD3, and SAC

The optimal pruning method for each baseline model was determined based on earning at least 90% average return of the corresponding baseline model while achieving the highest possible pruning percentage. The results of pruning experiments are presented in Table 2 and summarized comprehensively in Table 3. In Figure 2 and Figure 3, we present the effect of $L_1$ and $L_2$ pruning on the inference speed, energy usage, and memory usage. In these figures, we scaled the data according to the baseline.

| | TRPO | PPO | DDPG | TD3 | SAC |
|---|---|---|---|---|---|
| HalfCheetah | $L_1$ 5% | $L_1$ 10% | $L_2$ 5% | $L_2$ 5% | $L_1$ 5% |
| HumanoidStandup | $L_2$ 55% | $L_1$ 5% | $L_1$ 70% | $L_1$ 70% | $L_2$ 70% |
| Ant | $L_2$ 70% | $L2$ 5% | $L_1$ 10% | $L_2$ 5% | $L_2$ 25% |
| Humanoid | $L_2$ 30% | $L_1$ 5% | $L_2$ 5% | $L_1$ 10% | $L_2$ 25% |
| Hopper | $L_2$ 40% | $L_2$ 30% | $L_1$ 5% | $L_2$ 5% | $L_2$ 20% |

Table 3: Best pruning method for each environment and each model.

---
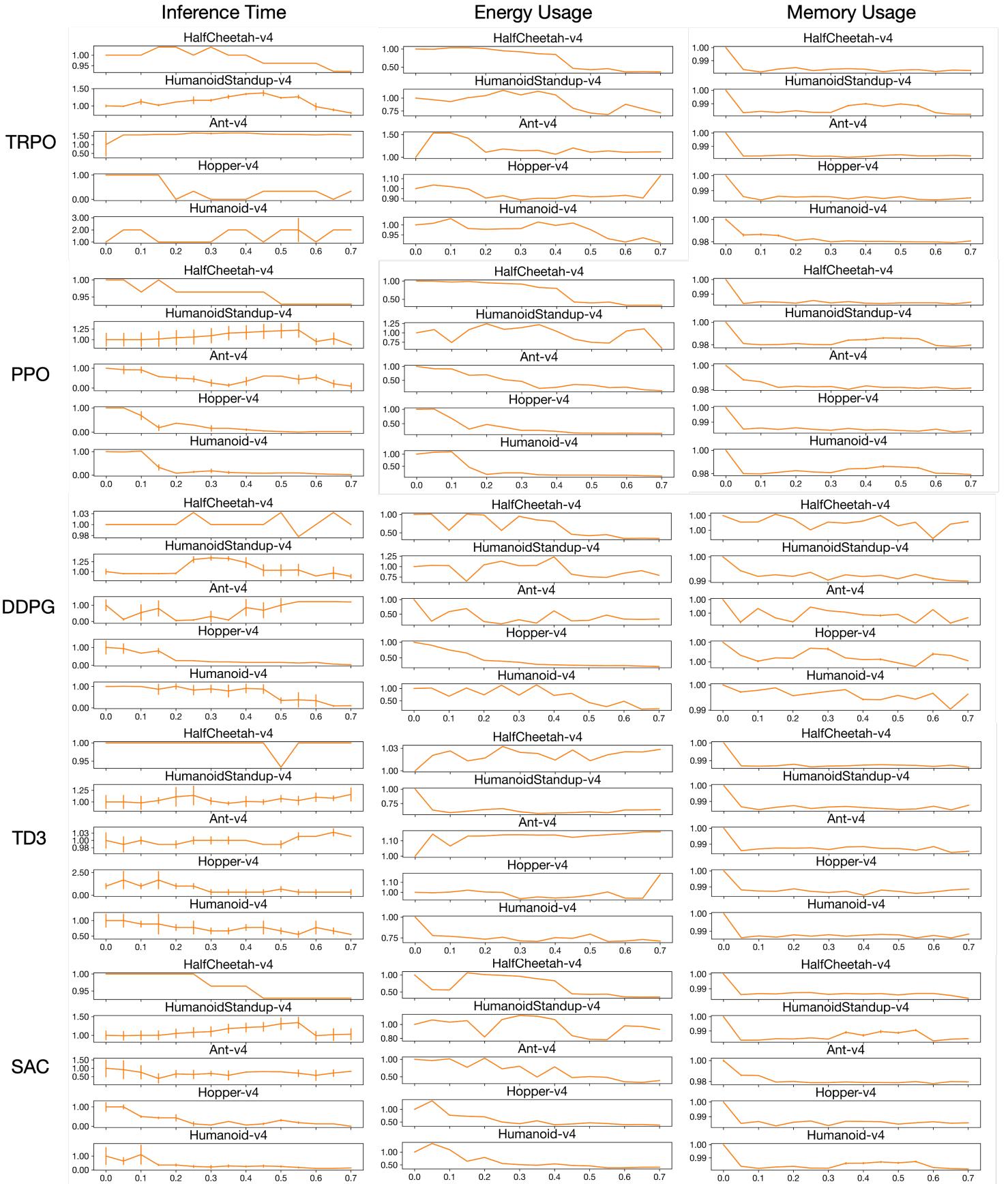
[8]https://github.com/VainF/Torch-Pruning

Figure 2: Inference time, energy usage and RAM of $L_1$ models, scaled by baseline models
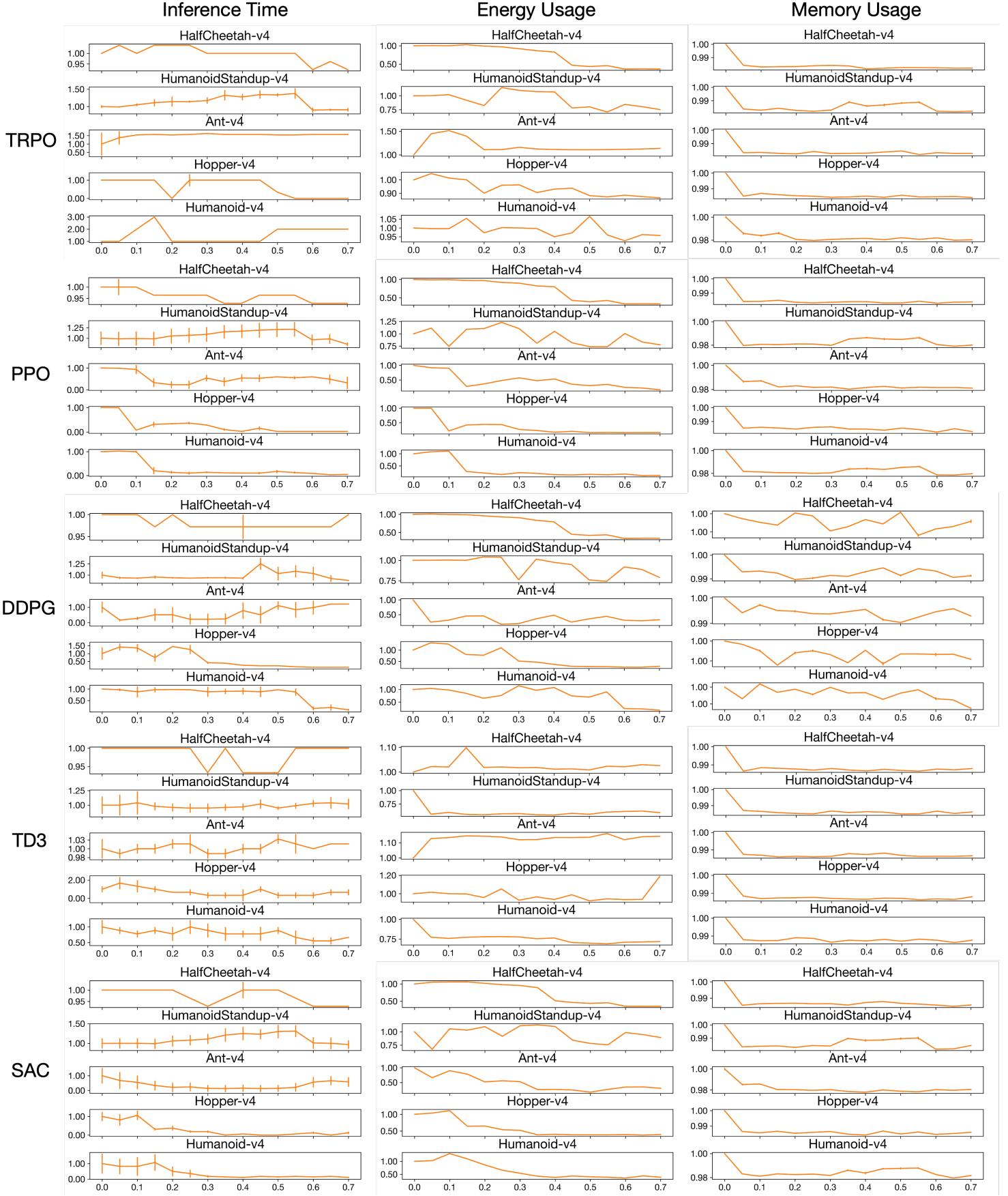
Figure 3: Inference time (in seconds), energy usage (in Joules) and memory utilization (in Megabytes) of $L_2$ models, scaled by baseline models.

# 4 Discussions and Findings

In this work, we studied two pruning approaches and three quantization approaches on five platforms (HalfCheetah-v4, HumanoidStandup-v4, Ant-v4, Hopper-v4, and Humanoid-v4) used for experimenting reinforcement learning methods and five common DRL methods (TRPO, PPO, DDPG, TD3). To our knowledge, this is the largest study performed on compressing DRL methods, and we listed our findings in this section. These findings could be used as a guideline for further studies that try to compress DRL methods.

*Pruning and quantization do not improve the energy efficiency and memory usage of DRL models.* While pruning and quantization reduce model size (see Table 3), they do not necessarily enhance the energy efficiency of DRL models due to the maintained or increased average return. Energy consumption tends to decrease only when there is a significant drop in average return, prompting the agent to terminate early and requiring less computation.

*Despite reducing model size, quantization does not improve memory usage, and pruning yields only a negligible 1% decrease in memory usage.* Results in Figure 1 present no changes in memory utilization in any platforms while applying quantization. Even PTDQ and PTSQ cause more memory utilization than the baseline method. This might be due to the overhead of the quantization library, and the way it is implemented is not optimized.

$L_2$ *pruning is favored over* $L_1$ *pruning for most of DRL models.* Table 2-3 illustrates that the optimal pruning method varies based on the DRL algorithm and environmental complexity. Most environments, except for those trained with SAC on HalfCheetah, allow for substantial pruning without a notable decline in average return, while PPO models exhibit lower pruning thresholds. In instances where $L_1$ pruning outperforms $L_2$, the average return values remain closely aligned. Generally, a 10% reduction in DRL model size through $L_2$ pruning is beneficial, although exceptions include PPO models applied to HalfCheetah environments.

*PTDQ emerges as the superior quantization method for DRL algorithms, whereas PTSQ is not recommended.* As shown in Table 1, 40% of our quantized models benefit from PTDQ, 36% from QAT, and only 24% from PTSQ. Our findings reveal that post-training dynamic quantization statistically outperforms the other methods, while post-training static quantization performs the worst, likely due to distribution shifts between the calibration data and the randomness that existed in RL environments.

*The Lottery ticket hypothesis [53] does not hold for DRL models.* The Lottery Ticket Hypothesis (LTH) in the context of neural networks suggests that within a large, randomly initialized network, there exists a smaller sub-network, typically around 10-20% of the original size, that, when trained in isolation, can achieve performance comparable to the original large network. This idea has significant implications for model quantization and pruning, two techniques used to reduce the size and computational requirements of neural networks. However, based on the results demonstrated in Table 2 demonstrate significant performance drops in most models after 50% pruning, contradicting the hypothesis's assertion that original network performance can persist even when pruned to less than 10%-20% of its original size. In particular, around 40% of the models don't survive after more than 5% pruning, but 80% of the models don't survive after 50%.

Our work has two limitations. First, by focusing on classical Mujoco environments with continuous action spaces, our work excludes discrete action spaces, which are common in video games or some decision-making scenarios. However, in these situations, task-specific methods might be employed for a satisfying performance, which adds additional complexities, and we might explore it in our future work. Moreover, we are limited to six simulated environments, this leaves an important aspect of real-world applicability unexplored. An ideal scenario is to experiment with this compression approach on a robot or drone in a real-world task and measure the differences in their performance.

# 5 Conclusion

In this paper, we examined the effect of quantization methods and pruning methods on deep reinforcement learning algorithms. While the effect depended on the specific DRL algorithm used and the environment in which the agent is trained, the results shared some common patterns. Quantization converts the floating point with 32 bits into an integer with 8 bits model and effectively shrinks the model size while maintaining acceptable performance. We found that PTDQ models generally had the best average return, while PTSQ models might suffer from distribution shifts and had poorer results. DepGraph pruned baseline models by constructing a dependency graph, and we experimented with $L_1$ and $L_2$ pruning. Experiments outlined that $L_2$ pruning was preferred for DRL algorithms on continuous action spaces, and in general, models benefited from 10% $L_2$ pruning with some exceptions. However, while pruning actually removed some neurons, it did not always result in inference speedup or energy saving.

# References

[1] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, L. Sifre, Dharshan Kumaran, Thore Graepel, Timothy P. Lillicrap, Karen Simonyan, and Demis Hassabis. Mastering chess and shogi

by self-play with a general reinforcement learning algorithm. *ArXiv*, abs/1712.01815, 2017.

[2] Peng Peng, Quan Yuan, Ying Wen, Yaodong Yang, Zhenkun Tang, Haitao Long, and Jun Wang. Multiagent bidirectionally-coordinated nets for learning to play starcraft combat games. *CoRR*, abs/1703.10069, 2017.

[3] Danijar Hafner, Timothy P. Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *CoRR*, abs/2010.02193, 2020.

[4] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155, 2022.

[5] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500, 2023.

[6] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.

[7] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2024.

[8] Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine*, 34(6):26–38, 2017.

[9] Marvin Zhang, Sharad Vikram, Laura M. Smith, P. Abbeel, Matthew J. Johnson, and Sergey Levine. Solar: Deep structured representations for model-based reinforcement learning. In *International Conference on Machine Learning*, 2018.

[10] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013.

[11] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA, 2020. Curran Associates Inc.

[12] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 1352–1361. JMLR.org, 2017.

[13] Richard S. Sutton, David A. McAllester, Satinder Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Neural Information Processing Systems*, 1999.

[14] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1889–1897, Lille, France, 07–09 Jul 2015. PMLR.

[15] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *ArXiv*, abs/1707.06347, 2017.

[16] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Tim Harley, Timothy P. Lillicrap, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, page 1928–1937. JMLR.org, 2016.

[17] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.

[18] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1582–1591. PMLR, 2018.

[19] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1856–1865. PMLR, 2018.

[20] Paul Francis Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *ArXiv*, abs/1706.03741, 2017.

[21] Alexandra Sasha Luccioni, Yacine Jernite, and Emma Strubell. Power hungry processing: Watts driving the cost of ai deployment? *arXiv preprint arXiv:2311.16863*, 2023.

[22] Alex de Vries. The growing energy footprint of artificial intelligence. *Joule*, 7(10):2191–2194, 2023.

[23] Peifeng Li, Jie Yang, Md Amirul Islam, and Suzhen Ren. Making ai less" thirsty": Uncovering and addressing the secret water footprint of ai models. *arXiv preprint arXiv:2304.03271*, 2023.

[24] Reza Rawassizadeh, Blaine Price, and Marian Petre. Wearables: Has the age of smartwatches finally arrived? *Communications of the ACM*, 58:45–47, 01 2015.

[25] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.

[26] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. In *International Conference on Learning Representations*, 2018.

[27] Zhenhua Liu, Yunhe Wang, Kai Han, Wei Zhang, Siwei Ma, and Wen Gao. Post-training quantization for vision transformer. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 28092–28103. Curran Associates, Inc., 2021.

[28] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[29] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.

[30] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. SmoothQuant: Accurate and efficient post-training quantization for large language models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 38087–38099. PMLR, 23–29 Jul 2023.

[31] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 10088–10115. Curran Associates, Inc., 2023.

[32] Srivatsan Krishnan, Maximilian Lam, Sharad Chitlangia, Zishen Wan, Gabriel Barth-Maron, Aleksandra Faust, and Vijay Janapa Reddi. Quarl: Quantization for fast and environmentally sustainable reinforcement learning. *Trans. Mach. Learn. Res.*, 2022, 2019.

[33] Je Yang, Seongmin Hong, and Joo-Young Kim. Fixar: A fixed-point deep reinforcement learning platform with quantization-aware training and adaptive parallelism. In *2021 58th ACM/IEEE Design Automation Conference (DAC)*, pages 259–264, 2021.

[34] Peng Hu, Xi Peng, Hongyuan Zhu, Mohamed M. Sabry Aly, and Jie Lin. Opq: Compressing deep neural networks with one-shot pruning-quantization. *ArXiv*, abs/2205.11141, 2021.

[35] Tailin Liang, John Glossner, Lei Wang, Shaobo Shi, and Xiaotong Zhang. Pruning and quantization for deep neural network acceleration: A survey. *Neurocomputing*, 461:370–403, 2021.

[36] Christos N. Mavridis and John S. Baras. Vector quantization for adaptive state aggregation in reinforcement learning. In *2021 American Control Conference (ACC)*, pages 2187–2192, 2021.

[37] Evans Miriti and Andrew Mwaura. Dynamic vector quantization for reinforcement learning (dvqrl). In *2018 5th International Conference on Soft Computing & Machine Intelligence (ISCMI)*, pages 38–42, 2018.

[38] Sherjil Ozair, Yazhe Li, Ali Razavi, Ioannis Antonoglou, Aaron Van Den Oord, and Oriol Vinyals. Vector quantized models for planning. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8302–8313. PMLR, 18–24 Jul 2021.

[39] Robert Dadashi, Léonard Hussenot, Damien Vincent, Sertan Girgin, Anton Raichuk, Matthieu Geist, and Olivier Pietquin. Continuous control with action quantization from demonstrations. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 4537–4557. PMLR, 17–23 Jul 2022.

[40] Yoonhee Gil, Jong-Hyeok Park, Jongchan Baek, and Soohee Han. Quantization-aware pruning criterion for industrial applications. *IEEE Transactions on Industrial Electronics*, 69(3):3203–3213, 2022.

[41] Dor Livne and Kobi Cohen. Pops: Policy pruning and shrinking for deep reinforcement learning. *IEEE Journal of Selected Topics in Signal Processing*, 14(4):789–801, 2020.

[42] Weiwei Zhang, Ming Ji, Haoran Yu, and Chenghui Zhen. Relp: Reinforcement learning pruning method based on prior knowledge. *Neural Processing Letters*, 55(4):4661–4678, 2023.

[43] Wensheng Su, Zhenni Li, Minrui Xu, Jiawen Kang, Dusit Tao Niyato, and Shengli Xie. Compressing deep reinforcement learning networks with a dynamic structured pruning method for autonomous driving. *ArXiv*, abs/2402.05146, 2024.

[44] Haoli Zhao, Jiqiang Wu, Zhenni Li, Wuhui Chen, and Zibin Zheng. Double sparse deep reinforcement learning via multilayer sparse coding and nonconvex regularized pruning. *IEEE Transactions on Cybernetics*, 53(2):765–778, 2023.

[45] Stephane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 627–635, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR.

[46] Hongjie Zhang, Zhuocheng He, and Jing Li. Accelerating the deep reinforcement learning with neural network compression. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2019.

[47] Mohammad Babaeizadeh, Iuri Frosio, Stephen Tyree, Jason Clemons, and Jan Kautz. GA3C: GPU-based A3C for deep reinforcement learning. *NIPS Workshop*, 2016.

[48] Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Variational dropout sparsifies deep neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, 2017.

[49] Jianping Gou, Baosheng Yu, Stephen Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129, 06 2021.

[50] Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards any structural pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16091–16101, June 2023.

[51] Mark Towers, Jordan K. Terry, Ariel Kwiatkowski, John U. Balis, Gianluca de Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Arjun KG, Markus Krimmel, Rodrigo Perez-Vicente, Andrea Pierré, Sander Schulhoff, Jun Jet Tai, Andrew Tan Jin Shen, and Omar G. Younis. Gymnasium, March 2023.

[52] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.

[53] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2018.