

Experiences with VITIS AI for Deep Reinforcement Learning

Nabayan Chaudhury
Department of CS
Virginia Tech
Blacksburg, VA, USA
nabayanc@vt.edu

Atharva Gondhalekar
Department of ECE
Virginia Tech
Blacksburg, VA, USA
atharval@vt.edu

Wu-chun Feng
Departments of CS and ECE
Virginia Tech
Blacksburg, VA, USA
wfeng@vt.edu

Abstract—Deep reinforcement learning has found use cases in many applications, such as natural language processing, self-driving cars, and spacecraft control applications. Many use cases of deep reinforcement learning seek to achieve inference with low latency and high accuracy. As such, this work articulates our experiences with the AMD Vitis AI toolchain to improve the latency and accuracy of inference in deep reinforcement learning.

In particular, we evaluate the soft actor-critic (SAC) model that is trained to solve the MuJoCo humanoid environment, where the objective of the humanoid agent is to learn a policy that allows it to stay in motion for as long as possible without falling over. During the training phase, we prune the model using the weight sparsity pruner from the Vitis AI optimizer at different timesteps. Our experimental results show that pruning leads to an improvement in the evaluation of the reinforcement learning policy, where the trained agent can remain balanced in the environment and accumulate higher rewards, compared to a trained agent without pruning. Specifically, we observe that pruning the network during training can deliver up to 20% better mean episode length and 23% higher reward (better accuracy), compared to a network without any pruning. Additionally, there is an improvement in decision-making latency up to 20%, which is the time between the observation of the agent’s state and a control decision.

Index Terms—reinforcement learning, humanoid, MuJoCo, network pruning, parallel computing, FPGA, GPU, Vitis AI

I. INTRODUCTION

In recent years, deep reinforcement learning (DRL) has found use cases in applications such as self-driving cars [1], natural language processing [2], and mission-critical tasks such as landing a spacecraft on celestial bodies [3]. Low-latency decision-making while maintaining the quality of solutions in resource and power-constrained environments is critical to the success of DRL algorithms in many applications, making fast and accurate DRL policy evaluation desirable.

Recent studies in the area of machine learning (ML) have explored the impact of pruning the trained network on the accuracy and performance of inference [4]–[6]. The process of pruning neural networks typically involves removing the connections between network layers, based on a pruning policy. Pruning has been shown to reduce the amount of computation necessary to generate the output of the inference,

The work detailed herein has been supported in part by NSF IUCRC CNS-1822080 via the NSF Center for Space, High-performance, and Resilient Computing (SHREC).

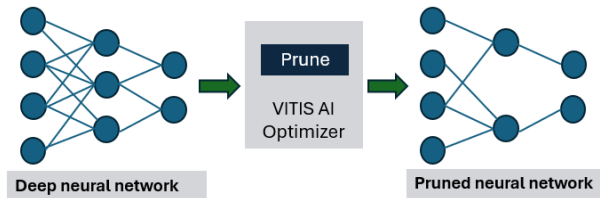


Fig. 1: Pruning a deep neural network using Vitis AI

and in many cases, maintain the accuracy of the inference [4], [5]. Fig. 1 shows an example of pruning a neural network using Vitis AI [7], a software stack developed by AMD for neural network inference on field-programmable gate arrays (FPGAs). In this work, we evaluate the efficacy of a weight sparsity pruner from the Vitis AI toolchain for the soft actor-critic (SAC) for DRL.

Specifically, we evaluate the soft actor-critic (SAC) model that is trained to solve the MuJoCo humanoid environment [8], where the humanoid agent learns a policy with the objective of remaining in a healthy state (i.e., remaining in motion) for as long as possible without falling over. We analyze the impact of pruning the network at various stages of the training phase and measure the performance of the pruned network using quantitative measures such as mean reward and mean episode length. We evaluate the performance impact of model pruning on an Nvidia RTX 3090 GPU, providing us insight into model pruning for DRL networks.

In all, we make the following contributions in this paper.

- Application of Vitis AI to prune DRL networks for higher reward and better mean episode length.
- Rigorous evaluation of the efficacy of pruning the trained model at various stages during the training phase.
- Up to 23% higher reward and 20% higher (i.e., better) mean episode length compared to the neural network without any pruning.
- Up to 20% faster decision-making compared to the neural network without any pruning.

II. RELATED WORK

We present related work in three parts: (1) reinforcement learning methods for the MuJoCo humanoid environment, (2)

prior work on exploring the effects of pruning in machine learning, and (3) existing studies that make use of Vitis AI.

A. Reinforcement Learning (RL) in MuJoCo Environment

Multi-joint dynamics with contact, or MuJoCo for short, is a general-purpose physics environment developed by Google-Deepmind [9]. This work focuses on the humanoid environment within MuJoCo, where the objective of the humanoid agent is to remain in a healthy state characterized by a continuous state of balance without falling over. The MuJoCo environments have been extensively used in machine-learning (ML) research. Wen et al. [10] use the MuJoCo environment to evaluate their multi-agent transformer network that casts multi-agent reinforcement learning (RL) into a sequence modeling problem. Shao et al. [11] present design-space exploration techniques for hardware acceleration of RL policy training and evaluate their hardware accelerator using MuJoCo environments. Liang et al. [12] present GPU-accelerated RL simulations using the MuJoCo environment.

B. Pruning in Machine Learning

Hoeffler et al. [4] survey many approaches to prune neural networks and explore multiple training strategies to achieve model sparsity. Chaturvedi et al. [5] explore the effects of introducing sparsity in a densenet and deconvolution network (DDNet). Obando-Ceron et al. [6] demonstrate that by removing network parameters during reinforcement learning (RL) policy training, it is possible to perform policy evaluations with better accuracy than evaluations with dense network counterparts. Inspired by the aforementioned studies, this work evaluates the impact of pruning DRL networks using Vitis AI.

C. Vitis AI

Vitis AI is a software stack developed by AMD for accelerating artificial intelligence inference on AMD FPGAs [7]. Ushiroyama et al. [13] use Vitis AI to implement convolutional neural network on FPGAs. Cabrera et al. [14] use Vitis AI for performing errant beam detection on AMD FPGAs.

Our work differs from these prior studies in the following ways. First, while Vitis AI is primarily focused on optimizing the trained networks for FPGAs, we show and evaluate our pipeline that integrates Vitis AI components with a GPU implementation. Second, while the effects of pruning, the use of Vitis AI, and DRL in the MuJoCo environment are individually well-explored subjects, we evaluate the multiplicative impact of pruning DRL networks for MuJoCo environments using Vitis AI.

III. BACKGROUND

A. Off-Policy Reinforcement Learning and Soft Actor-Critic

Reinforcement learning (RL) is a branch of machine learning (ML), where an agent learns to make optimal decisions by interacting with an environment and with the goal of maximizing a cumulative *reward* [15]. Deep reinforcement learning (DRL) combines RL with deep learning (DL), allowing deep neural networks to learn the policy from instances of the

environment without constructing a complete state space of the environment. The general training mechanism in RL involves the agent making decisions in the environment, receiving rewards, and updating its policy to maximize cumulative rewards over time. Formally, the reinforcement learning (RL) problem can be defined as a policy search in a Markov decision process (MDP) defined by a tuple (S, A, P, R, γ) , where S is the set of states, A is the set of actions, P represents the state transition probabilities, R is the reward function, and γ is the discount factor. The objective of the agent is to learn a policy $\pi(a | s)$ that maximizes the expected cumulative reward, defined as the return G_t :

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

where R_{t+k+1} is the reward received $k + 1$ steps after time t , and $\gamma \in [0, 1)$ is the discount factor that determines the importance of future rewards. The goal is to find the optimal policy π^* that maximizes the expected return from any initial state s_0 :

$$\pi^* = \arg \max_{\pi} \mathbb{E}[G_t | \pi] \quad (2)$$

State- and action-value functions estimate the expected return, and these three sets of equations govern the learning objective. The state-value function $V_{\pi}(s)$ and the action-value function $Q_{\pi}(s, a)$ are defined as:

$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s] \quad (3)$$

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a] \quad (4)$$

As the agent spends time performing actions in the environment, it updates its policy and value estimates iteratively as a coupled optimization objective.

RL algorithms can be designed in several ways. *Online reinforcement learning* involves continuous interaction of the agent with the environment, which allows the agent to learn and update its policy in real time. This enables the agent to adapt quickly to environment changes. Examples include Q-learning (DQN) [16] and actor-critic algorithms [17], [18]. *Offline reinforcement learning*, in contrast, utilizes pre-collected data to train the agent, where the agent has no further interactions with the environment. Offline algorithms can run into sub-optimal policy performance when the agent encounters novel states [19]. Notable offline algorithms include conservative Q-learning (CQL) [20] and batch-constrained deep Q-learning (BCQN) [21].

In addition, algorithms can be either *on-policy* or *off-policy*. On-policy algorithms, like proximal policy optimization (PPO) [22] update the policy based on actions taken by the current policy. These methods require consistent interaction with the environment to get accurate policy updates, but they have a tendency to explore the environment less as they rely on the current policy's knowledge of the environment. In contrast, off-policy algorithms learn from actions outside the current policy, encouraging exploration and improving

sample efficiency as the agent can learn from a diverse set of actions [18], [23], [24].

Our choice of algorithm is motivated by our interest in real-world, physics-based scenarios with large action spaces and complex environments, for which online, off-policy, and model-free algorithms are particularly advantageous. This led to the selection of soft actor critic (SAC) [18] for our experiments. SAC maximizes a tradeoff between the expected reward and the entropy of the policy, which promotes exploration and prevents premature convergence to suboptimal policies.

B. Simulated Physics Environments

Once the RL algorithms have been trained, it is essential to have a method for benchmarking them in a controlled, reproducible setting without interacting with the real world. Such environments are generally dynamical systems with complex control laws, like *rigid body* dynamical systems (objects that do not deform and interact through collisions and forces) and *articulated body* dynamical systems (interconnected rigid bodies that exhibit human or animal motion). These systems have their own *action space* (which is a superset of the policy action space A) and *observation space* (which is a superset of the policy state space S), where:

- Action space of the system (A_S): set of all possible actions an agent can take in the environment.
- State space of the system (S_S): set of all possible configurations of the environment.

Continuous action spaces consist of a range of values in a given interval as control inputs and closely emulate real-world, physics-based application scenarios. For example, the Multi-joint dynamics with contact (MuJoCo) [25] engine provides several such articulated rigid body continuous control environments that are of interest to this study. *Discrete* action spaces consist of a finite set of actions, where each action is represented as a single binary choice, making the space lower-dimensional. This simplifies the decision-making process but may not capture the complexity of actions required in a more sophisticated environment [26].

C. Vitis AI

Fig. 2 describes Vitis AI, a comprehensive development stack developed by AMD for accelerating artificial intelligence (AI) inference on AMD hardware, including FPGAs (field-programmable gate arrays) and SoCs (system on a chip) [7]. It is designed with efficiency and acceleration in mind, allowing users to deploy accelerated machine-learning (ML) models on AMD FPGAs. It provides an end-to-end solution — from model optimization to deployment — where users can achieve high performance, low latency, and efficient resource utilization. Fig. 2 shows the following key components of Vitis AI (or VAI for short):

- *VAI Model Zoo*: A set of pre-trained and optimized models that can be easily deployed on AMD FPGA devices.
- *VAI Optimizer*: A tool capable of performing pruning on AI models, reducing model size, and improving inference

speed without affecting accuracy. This paper focuses on this particular tool.

- *VAI Quantizer*: A tool that quantizes AI models by reducing precision from `fp32` (single-precision floating point) to lower precision formats like `int8`, accelerating computation and reducing memory usage.
- *VAI Compiler*: A tool that compiles the reduced model into a deployable model supported by AMD FPGA hardware, e.g., AMD Deep-Learning Processing Unit (DPU) [7].
- *VAI Profiler*: An application-level tool that allows for thorough profiling of deployable models for inference, giving insights into further optimizations for maximum performance.
- *VAI Library*: A library of APIs for easy utilization of the VAI software stack. It supports models in the Model Zoo as well as custom models.
- *VAI Runtime*: A runtime environment that executes compiled models on AMD FPGA hardware and capable of efficient resource allocation and task scheduling.
- *Deep-Learning Processing Units (DPUs)*: General-purpose AI inference accelerators that target convolution neural networks (CNNs), allowing for simultaneous deployment and inference.

To date, Vitis AI has been used to implement convolutional neural networks [29], accelerate existing object detection algorithms [30], and perform embedded object detection with custom model architectures [31]. In all cases, the reported results highlight a significant decrease in power consumption and increase in throughput. Additionally, model quantization has been shown to increase robustness against adversarial examples [32], and the Vitis AI development stack has been used in real-world errant beam detection, showing fast performance and high accuracy [33].

IV. CASE STUDY: MUJOCo HUMANOID

In this section, we use the “multi-joint dynamics with contact” environment (i.e., MuJoCo) as our case study to illustrate the efficacy of Vitis AI.

A. Environment

The MuJoCo humanoid environment simulates a 3D bipedal humanoid robot designed to mimic a human standing upright. [34]. Fig. 3 shows the humanoid agent in different states that it can end up in through an episode of training [8]. The robot consists of a torso with a pair of legs and arms. Each leg has three segments (thigh, shin and foot), and each arm has two segments (upper and lower arm).

The primary goal in this environment is to control the humanoid to stay upright by balancing itself as long as possible without toppling over. The following summarizes the action and state spaces:

- **Action Space:**
 - *Type*: Continuous
 - *Dimension (Degrees of freedom)*: 17
 - *Range*: Box(-0.4, 0.4, (17,)), float32)

Vitis AI: Unified AI Inference Solution Stack

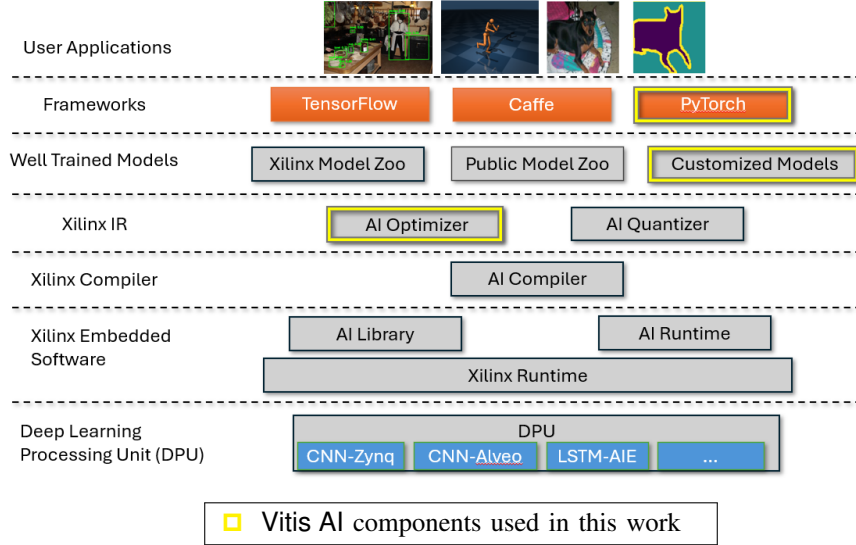


Fig. 2: Vitis AI: Unified AI inference solution stack [7], [27], [28]

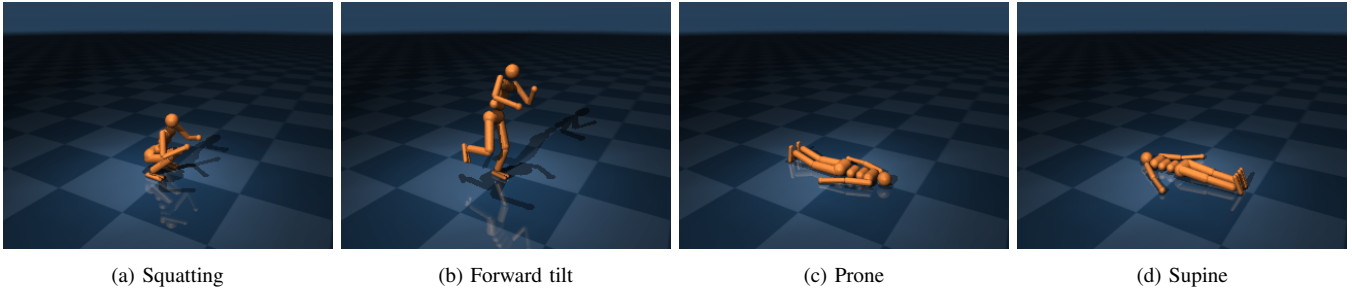


Fig. 3: MuJoCo humanoid exhibiting different poses [8], [9]

- *Description*: The action space is defined by the torques acting at each of the humanoid’s joints

- **Observation Space:**

- *Type*: Continuous
- *Dimension*: 376
- *Range*: Box(-inf, inf, (376,), float64)
- *Description*: The observation space consists of the position and velocity values of the various joints and body parts. These in turn define the state the humanoid is in at any given time.

The reward structure in the humanoid encourages stable and efficient balancing. It is divided into a *healthy reward* that is a fixed reward for every timestep that the humanoid remains standing, a *forward reward* that is calculated based on the forward displacement of the humanoid’s center of mass and encourages forward movement, a *control cost* penalty for using excessive control forces, and a *contact cost* penalty for high external forces. This naturally translates to balancing being human-like and not forced by excessive inputs. Each episode terminates when either the humanoid reaches 1,000

timesteps, or the humanoid becomes *unhealthy*, signified by the z-coordinates of the torso falling outside a healthy range. Because the goal of the agent is to stay upright as long as possible, performance evaluation can be done by measuring the *average episode lengths* and *average cumulative rewards*. If the episodes run for longer timesteps, the cumulative reward increases and the agent is capable of staying upright longer.

B. Algorithm

For the purpose of our experiments, we use the Stable-Baselines 3 (SB3) library and its implementation of the Soft Actor-Critic (SAC) algorithm [18], [35]. SB3 provides a comprehensive set of algorithm implementations for fast development and deployment, while supporting custom policy networks and in-house benchmarking and evaluation. We choose SAC because of its model-free, off-policy, and online nature, making it ideal for solving large, physics-based environments that require exploration to reach an optimal policy. We train the agent using SAC for a total of 100,000 timesteps in each experiment, pruning the model after we reach 10%, 20%, 50%, 75% and 100% of training timesteps. Pruning is

done using the Sparse Pruner from the Vitis AI Optimizer, setting the weight sparsity to 0.5. The complete details of the algorithm can be found in Algorithm 1.

Algorithm 1: Training and pruning SAC model on MuJoCo humanoid with Vitis AI

Input : Total timesteps `TOTAL_TIMESTEPS`,
Pruning timestep `PRUNING_TIMESTEP`,
Number of evaluation episodes
`NUM_EVAL_EPISODES`

Output: Trained and pruned SAC model, Evaluation results

Data: Training and evaluation environments, SAC model, Vitis AI optimizer, Callbacks

- 1 **Initialize:** Training and evaluation environments, action noise, SAC model, and callbacks;
- 2 **for** `timestep` \leftarrow 0 **to** `TOTAL_TIMESTEPS - 1` **do**
- 3 **if** `timestep` == `PRUNING_TIMESTEP` **then**
- 4 Prune the model using Vitis AI Optimizer with specified sparsity;
- 5 **end**
- 6 **end**
- 7 **for** `episode` \leftarrow 0 **to** `NUM_EVAL_EPISODES - 1` **do**
- 8 Evaluate the trained model for one episode;
- 9 Log episode results (rewards, episode lengths);
- 10 **end**

TABLE I: Summary of training and inference performance for humanoid environment

Model (SAC policy network)	Training time (s)	Evaluation time (s)	Mean cumulative reward	Mean episode length	Mean inference latency (ms)
Baseline (no pruning)	890.60	11.78	644.99	131.07	0.708
Pruned after 10000 (10%) timesteps	1018.86	13.28	646.67	128.85	0.570
Pruned after 20000 (20%) timesteps	1038.79	14.19	781.34	139.47	0.572
Pruned after 50000 (50%) timesteps	978.01	16.72	812.62	159.68	0.569
Pruned after 75000 (75%) timesteps	946.83	13.75	662.24	131.31	0.567
Pruned after 100000 timesteps (full training)	894.58	14.18	696.99	134.54	0.563

Mean reward and mean episode length: higher is better
Training and evaluation time measured on NVIDIA RTX 3090 GPU.

V. EVALUATION

We evaluate the performance of six different SAC models trained on the MuJoCo Humanoid-v4 environment on an

NVIDIA RTX 3090 GPU by running each policy on 100 unique evaluation environments to ensure robustness of results. Our trained models are in PyTorch, and we use 100 unique instances of the MuJoCo Humanoid for each model as our evaluation environments.

The models under evaluation are as follows:

- Baseline (no pruning)
- Policy network pruned after 10% training timesteps
- Policy network pruned after 20% training timesteps
- Policy network pruned after 50% training timesteps
- Policy network pruned after 75% training timesteps
- Policy network pruned after 100% training timesteps, i.e., after the model has trained to completion.

Model performance is measured by the *mean cumulative reward* and the *mean episode length*. Better policy evaluation is indicated by larger mean rewards and longer episode lengths as the agent can stay in a ‘healthy state’ for an extended period, where it can remain in motion without falling over. Fig. 6 shows the accumulated reward evolution over time during the training process for the baselines compared to models pruned at a specific timestep. Additionally, we define a metric *mean inference latency*, which measures the time delay between the observation of the agent’s current state and the policy decision being made. As such, lower latency suggests that the policy network is quicker at taking decisions.

These trained policies are each run on 100 evaluation environment, with the values for mean accumulated rewards and mean episode lengths reported in Table I. Note that the training times increase due to the pruning process itself requiring some overhead, and the evaluation time increases with improved performance as the agent under optimal policy now spends more time in the environment and accumulates higher rewards. Fig. 4a and 4b show the distributions of rewards and episode lengths of the 100 evaluation instances. In particular, pruning the model halfway through the training process leads to a model with the best evaluation performance overall, suggesting a trade-off point where the pre-trained, non-pruned model after pruning can generalize quickly with sufficient further training.

There is also noticeable reduction of decision-making latency by $\approx 20\%$ when comparing the pruned models to the baseline architecture, as shown in 4c. This suggests that a pruned network, while simpler in architecture, delivers faster decision-making without any trade-offs in evaluation performance. Fig. 5 shows the agent trying to stay upright. At $T=16.72s$, when the episode ends, the agent reaches a state where it can no longer correct itself to keep upright.

Overall, the results indicate that model pruning generally improves the evaluation performance of SAC models in the MuJoCo Humanoid environment and that there exists an optimal pruning timestep where performance gains can be maximized.

VI. FUTURE WORK

We plan to extend our investigation into the other steps of the Vitis AI toolchain, specifically the quantizer and the compiler. We want to apply the complete Vitis AI workflow to

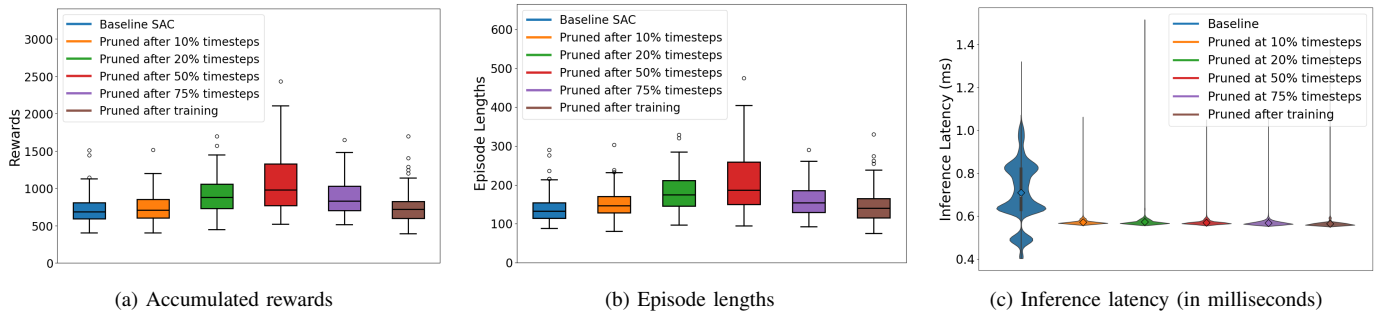


Fig. 4: Policy evaluations for networks pruned at different timesteps

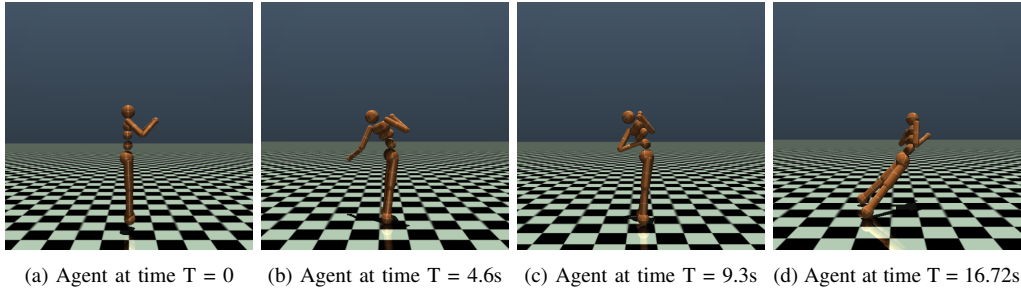


Fig. 5: The evolution of the agent over time for the best performing policy network (pruned at 50% training time). The episode ends when the agent reaches a state where it can no longer correct itself to stay upright

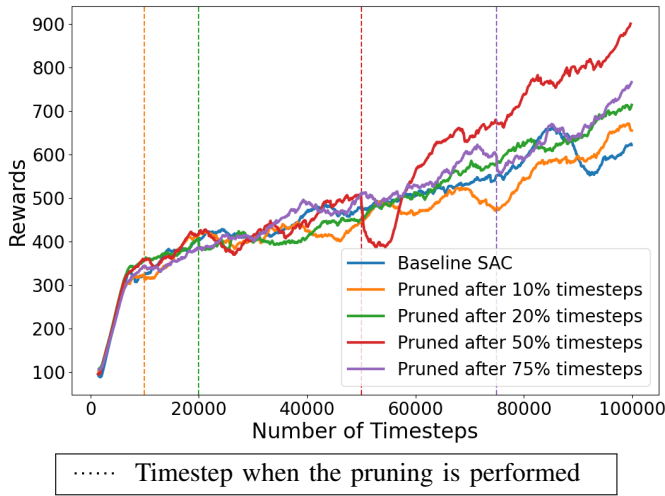


Fig. 6: Training reward evolution over time

several off-policy, online deep reinforcement learning (DRL) algorithms.

VII. CONCLUSION

In this study, we explore the impact of model pruning on the performance of Soft Actor-Critic (SAC) algorithms trained on the MuJoCo Humanoid environment. We employ the Sparse Pruner Optimizer from the Vitis AI toolchain at various stages of the training process to investigate the effect of pruning on accumulated rewards and episode lengths and evaluate these models on 100 unique evaluation environments each

to ensure robustness. All models were trained and evaluated on an Nvidia RTX 3090 GPU. Our findings indicate that pruning generally enhances the performance of SAC models in terms of overall reward as well as decision making speed. Specifically, pruning at 50% of the total timesteps resulted in a 23% increase in accumulated rewards and 20% increase in episode lengths, which is the maximum observed performance improvement, suggesting that this interval may strike a favorable balance between sufficient training and effective pruning. We also notice that pruning overall reduces the inference latency (between observation and action) for the learned policy, speeding up decision making by 20%. However, it is important to note that this observation is based on empirical evidence from a limited set of pruning intervals.

Through this study, we aim to provide a foundation for further research into the Vitis AI toolchain for deep reinforcement learning inference. In future work, we plan to incorporate other steps in the toolchain: model quantization, compilation and deployment to FPGA hardware for inference, and compare performance to traditional CPU and GPU implementations. Furthermore, we plan to include other state-of-the-art deep RL algorithms in our analysis, such as Proximal Policy Optimization (PPO), Advantage Actor Critic (A2C) and Twin-Delayed DDPG (TD3).

ACKNOWLEDGEMENT

The work was supported in part by NSF I/UCRC CNS-1822080 via the NSF Center for Space, High-performance, and Resilient Computing (SHREC).

REFERENCES

- [1] A. R. Fayjie, S. Hossain, D. Oualid, and D.-J. Lee, "Driverless Car: Autonomous Driving Using Deep Reinforcement Learning in Urban Environment," in *2018 15th International Conference on Ubiquitous Robots (UR)*, 2018, pp. 896–901.
- [2] V. Uc-Cetina, N. Navarro-Guerrero, A. Martin-Gonzalez, C. Weber, and S. Wermter, "Survey on reinforcement learning for language processing," *Artificial Intelligence Review*, vol. 56, no. 2, pp. 1543–1575, Feb 2023. [Online]. Available: <https://doi.org/10.1007/s10462-022-10205-5>
- [3] S. Gadgil, Y. Xin, and C. Xu, "Solving The Lunar Lander Problem under Uncertainty using Reinforcement Learning."
- [4] T. Hoefler, D. Alistarh, T. Ben-Nun, N. Dryden, and A. Peste, "Sparsity in deep learning: pruning and growth for efficient inference and training in neural networks," *J. Mach. Learn. Res.*, vol. 22, no. 1, Jan 2021.
- [5] A. Chaturvedi, G. Cao, and W. chun Feng, "Optimizing Deep Learning for Biomedical Imagin," in *International Conference on Computational Advances in Bio and medical Sciences*, December 2023.
- [6] J. Obando-Ceron, A. Courville, and P. S. Castro, "In value-based deep reinforcement learning, a pruned network is a good network," 2024. [Online]. Available: <https://arxiv.org/abs/2402.12479>
- [7] (2024) Vitis AI. Advanced Mirco Devices (AMD). [Online]. Available: <https://github.com/Xilinx/Vitis-AI>
- [8] E. Todorov, T. Erez, and Y. Tassa, "MuJoCo: A physics engine for model-based control," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 5026–5033.
- [9] Google-Deepmind. (2024) MuJoCo. [Online]. Available: <https://github.com/google-deepmind/mujoco>
- [10] M. Wen, J. Kuba, R. Lin, W. Zhang, Y. Wen, J. Wang, and Y. Yang, "Multi-Agent Reinforcement Learning is a Sequence Modeling Problem," in *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35. Curran Associates, Inc., 2022, pp. 16509–16521. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2022/file/69413f87e5a34897cd010ca698097d0a-Paper-Conference.pdf
- [11] S. Shao, J. Tsai, M. Mysior, W. Luk, T. Chau, A. Warren, and B. Jeppesen, "Towards Hardware Accelerated Reinforcement Learning for Application-Specific Robotic Control," in *2018 IEEE 29th International Conference on Application-specific Systems, Architectures and Processors (ASAP)*, 2018, pp. 1–8.
- [12] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox, "GPU-Accelerated Robotic Simulation for Distributed Reinforcement Learning," in *Conference on Robot Learning*, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:53084610>
- [13] A. Ushiroyama, M. Watanabe, N. Watanabe, and A. Nagoya, "Convolutional neural network implementations using vitis ai," in *2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC)*, 2022, pp. 0365–0371.
- [14] A. M. Cabrera, Y. A. Yucasan, F. Y. Liu, W. Blokland, and J. S. Vetter, "Errant Beam Detection Using the AMD Versal ACAP and Vitis AI," in *2023 IEEE High Performance Extreme Computing Conference (HPEC)*, 2023, pp. 1–6.
- [15] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. The MIT Press, 2018. [Online]. Available: <http://incompleteideas.net/book/the-book-2nd.html>
- [16] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb 2015. [Online]. Available: <https://doi.org/10.1038/nature14236>
- [17] V. Konda and J. Tsitsiklis, "Actor-critic algorithms," in *Advances in Neural Information Processing Systems*, S.olla, T. Leen, and K. Müller, Eds., vol. 12. MIT Press, 1999. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/1999/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf
- [18] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *ICML*, ser. Proceedings of Machine Learning Research, J. G. Dy and A. Krause, Eds., vol. 80. PMLR, 2018, pp. 1856–1865. [Online]. Available: <http://dblp.uni-trier.de/db/conf/icml/icml2018.html#HaarnojaZAL18>
- [19] S. Levine, A. Kumar, G. Tucker, and J. Fu, "Offline reinforcement learning: Tutorial, review, and perspectives on open problems," *ArXiv*, vol. abs/2005.01643, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:218486979>
- [20] A. Kumar, A. Zhou, G. Tucker, and S. Levine, "Conservative q-learning for offline reinforcement learning," in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, ser. NIPS '20. Red Hook, NY, USA: Curran Associates Inc., 2020.
- [21] S. Fujimoto, D. Meger, and D. Precup, "Off-policy deep reinforcement learning without exploration," in *International Conference on Machine Learning*, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:54457299>
- [22] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms." *CoRR*, vol. abs/1707.06347, 2017. [Online]. Available: <http://dblp.uni-trier.de/db/journals/corr/corr1707.html#SchulmanWDRK17>
- [23] 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." in *ICLR*, Y. Bengio and Y. LeCun, Eds., 2016. [Online]. Available: <http://dblp.uni-trier.de/db/conf/iclr/iclr2016.html#LillicrapHPHETS15>
- [24] S. Dankwa and W. Zheng, "Twin-delayed ddpq: A deep reinforcement learning technique to model a continuous movement of an intelligent robot agent," in *Proceedings of the 3rd International Conference on Vision, Image and Signal Processing*, ser. ICVISP 2019. New York, NY, USA: Association for Computing Machinery, 2020. [Online]. Available: <https://doi.org/10.1145/3387168.3387199>
- [25] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012, pp. 5026–5033.
- [26] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," 2016, cite arxiv:1606.01540. [Online]. Available: <http://arxiv.org/abs/1606.01540>
- [27] G. Jocher, A. Stoken, J. Borovec, NanoCode012, ChristopherSTAN, L. Changyu, Laughing, tkianai, A. Hogan, lorenzomammama, yxNONG, AlexWang1900, L. Diaconu, Marc, wanghaoyang0106, ml5ah, Doug, F. Ingham, Frederik, Guilhen, Hatovix, J. Poznanski, J. Fang, L. Y. , changyu98, M. Wang, N. Gupta, O. Akhtar, PetrDvoracek, and P. Rai, "ultralytics/yolov5: v3.1 - Bug Fixes and Performance Improvements," Oct. 2020. [Online]. Available: <https://doi.org/10.5281/zenodo.4154370>
- [28] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, "TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems," 2015, software available from tensorflow.org. [Online]. Available: <https://www.tensorflow.org/>
- [29] A. Ushiroyama, M. Watanabe, N. Watanabe, and A. Nagoya, "Convolutional neural network implementations using vitis ai," in *2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC)*, 2022, pp. 0365–0371.
- [30] J. Wang and S. Gu, "FPGA Implementation of Object Detection Accelerator Based on Vitis-AI," in *2021 11th International Conference on Information Science and Technology (ICIST)*, 2021, pp. 571–577.
- [31] M. Machura, M. Danilowicz, and T. Kryjak, "Embedded object detection with custom littenet, finn and vitis ai dnn accelerators," *Journal of Low Power Electronics and Applications*, vol. 12, no. 2, 2022. [Online]. Available: <https://www.mdpi.com/2079-9268/12/2/30>
- [32] Y. Fukuda, K. Yoshida, and T. Fujino, "Evaluation of model quantization method on vitis-ai for mitigating adversarial examples," *IEEE Access*, vol. 11, pp. 87200–87209, 2023.
- [33] A. M. Cabrera, Y. A. Yucasan, F. Y. Liu, W. Blokland, and J. S. Vetter, "Errant beam detection using the amd versal acap and vitis ai," in *2023 IEEE High Performance Extreme Computing Conference (HPEC)*, 2023, pp. 1–6.
- [34] Y. Tassa, T. Erez, and E. Todorov, "Synthesis and stabilization of complex behaviors through online trajectory optimization," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012, pp. 4906–4913.
- [35] 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>