
FPGA-Gym: An FPGA-Accelerated Reinforcement Learning Environment Simulation Framework

Jiayi Li[◇], Hongxiao Zhao[◇], Wenshuo Yue[◇], Yihan Fu[◇], Daijing Shi[◇], Anjunyi Fan[◇],
Qinghao Wang[◇], Yaodong Yang^{◇†}, Bonan Yan^{◇*}

[◇]Peking University [†]Beijing Institute for General AI

*Corresponding author: bonanyan@pku.edu.cn

Abstract

1 Reinforcement learning (RL) faces the key challenge of balancing exploration (gather-
2 ing information about the environment through trials) and exploitation (exploiting
3 current knowledge to maximize rewards), especially in the open world. To boost
4 exploration efficiency, parallel environment execution is a widely used technique
5 that instantiates multiple environments and executes them in parallel. However,
6 this is computationally challenging using CPUs limited by the total thread numbers.
7 In this work, we present FPGA-Gym, an FPGA-CPU joint acceleration framework
8 for accelerating RL environment parallel executions. By offloading environment
9 steps to FPGA hardware, FPGA-Gym achieves a 4.36 to 972.6× speedup over
10 the existing fastest software-based framework for parallel environment execution.
11 Moreover, FPGA-Gym is a general RL acceleration framework compatible with
12 existing RL algorithms and frameworks. Its modular and parameterized features
13 allow users to conveniently customize new environments without extensive FPGA
14 knowledge. We demonstrate multiple representative RL benchmarks (e.g. Cartpole,
15 CliffWalking, Seaquest etc.) with Deep Q-Network and proximal policy optimiza-
16 tion algorithms. Additionally, we provide a standard environment library similar to
17 Gymnasium based on FPGA-Gym framework. The framework and the library are
18 available at https://github.com/Selinaee/FPGA_Gym.

19 1 Introduction

20 Reinforcement learning (RL) is highly effective in addressing sequential decision-making problems. It
21 empowers artificial agents to make decisions within an environment to maximize cumulative rewards
22 over time. RL faces the key challenge of balancing exploration (gathering information about the
23 environment through trials) and exploitation (exploiting current knowledge to maximize rewards) [1,
24 2], especially in the open world. To boost exploration efficiency, parallel environment execution is a
25 widely used technique that instantiates multiple environments and executes them in parallel, allowing
26 multiple agents to interact with the environment simultaneously, thereby rapidly accumulating
27 experience (Fig. 1(a)). However, updating a large number of environments in parallel places a
28 substantial demand on computational resources [3–5]. Fig. 1(b) reveals that parallel environment step
29 computation inflicts long latency, highlighting the significant computational demands required for
30 environment evolution.

31 The key challenge in parallel RL environment execution lies in rapidly updating numerous parallel
32 environments [3, 4, 6, 7]. On conventional hardware platforms, multithreading is the major tackle
33 for parallelism. Central processing units (CPUs) can only manage a certain level of parallelism
34 with limited thread numbers. Graphics processing units (GPUs) [8] and tensor processing units
35 (TPUs) are also under intense investigation to parallelize environment execution for RL. However,
36 RL environment execution presents unique challenges in sequential Markov decision process (MDP)

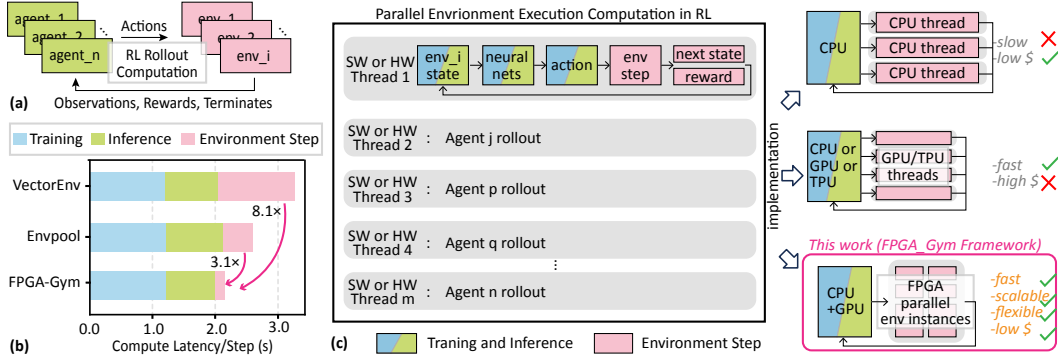


Figure 1: (a) Rollout computation of RL; (b) Compute latency breakdown in different RL acceleration frameworks; (c) Illustration of typical realizations of parallel RL environment execution.

37 computation [9–12]. The agent-environment interactions are modeled as Markov chains, characterized
 38 by iterative dynamic state transitions, such as environment resets and reward functions [13–15]. These
 39 elements are difficult to parallelize and scale efficiently on GPUs or TPUs (Fig. 1c).

40 To address the challenge, this work proposes to use field-programmable gate arrays (FPGAs) as the
 41 hardware acceleration foundation [16–21]. FPGAs are programmable semiconductor chips consisting
 42 of reconfigurable digital logic circuit blocks and interconnects. **To the best of our knowledge, this
 43 work presents the first holistic FPGA-CPU joint acceleration framework, termed as FPGA-Gym,
 44 for massively instantiating and parallelizing RL environment computation.** It offloads the
 45 execution of environment steps onto FPGAs via a high-bandwidth peripheral component interconnect
 46 express (PCIe) interface to the CPU. This work aims to overcome the limitations of CPUs and GPUs
 47 in handling logically complex and data-intensive RL tasks.

48 FPGA-Gym is fast and scalable. We propose pipeline and time division multiplexing techniques for
 49 acceleration [22]. Compared to EnvPool, the fastest CPU-based environment parallelism library
 50 at present, FPGA-Gym achieves a 4.36 to 972.6 \times speedup. In a grid-world task like *CliffWalking*,
 51 FPGA-Gym can instantiate more than 8,000 parallel environments in one FPGA chip.

52 FPGA-Gym is easy to customize. We implement FPGA-Gym in parameterized and modular manners. It
 53 provides a plug-in usage for users without FPGA knowledge. Customizing new environments on top
 54 of FPGA-Gym is convenient, only requiring revising the computing module of a single environment
 55 inside our provided template. Other parameterized parallel operations and interactions between FPGA
 56 and CPU are already defined in the framework backbone code.

57 FPGA-Gym is general and compatible with the existing algorithms and libraries. We release an
 58 open-source FPGA-Gym library, including various types of RL environments. Moreover, FPGA-Gym
 59 focuses on the environment execution and is compatible to the existing learning algorithms, e.g. deep
 60 Q-network (DQN) [23, 24] and proximal policy optimization (PPO) [25, 26], without modifying the
 61 learning framework source code.

62 2 Related Works

63 **CPU-based RL frameworks:** Gymnasium [27], VectorEnv and Sample Factory [28] offer
 64 multithreading or multiprocessing to run each environment in an independent thread or process.
 65 However, running programs with tens or hundreds of threads is troublesome in that the users have to
 66 take care of thread synchronization, context-switching overhead, and memory bandwidth limitations.
 67 EnvPool [22] leverages C++ backend for optimized parallel computing and memory management,
 68 achieving a 2.8 \times speedup compared to its Python counterpart. Nonetheless, it requires the developer
 69 to manually translate the Python-written environment into C++. As the number of instantiated
 70 environments increases, the parallelization becomes increasingly constrained by the fixed number of
 71 CPU threads, thereby limiting its scalability.

72 **GPU-/TPU- based RL frameworks:** GPUs/TPUs demonstrate significant performance boosts in
 73 computing large-scale linear algebra. Naturally, researchers would like to apply this advantage to RL
 74 acceleration. Researchers have adapted the source code of several pure-compute RL environments,
 75 e.g. *MuJoCo*, etc. [8, 29–42], to be compatible with GPU/TPU using CUDA or JAX toolkits for
 76 parallelism. Complex real-world or intricate logical tasks can be challenging to simulate using
 77 GPUs or TPUs due to their inefficiency in parallelizing sequential logic operators [43]. Moreover,
 78 the process of deploying environments to these hardware accelerators involves significant software
 79 engineering efforts, including the need to work with different programming models and low-level
 80 libraries. These challenges and limitations can make it difficult to fully leverage the computational
 81 power of GPUs and TPUs for such purposes.

82 3 FPGA-Gym Framework

83 3.1 User workflow

84 FPGA-Gym accommodates the needs of two user groups:

85 **Standard environment users**, who would like to accelerate the RL rollouts and training, can
 86 simply ① plug the FPGA board into a computer motherboard via PCIe interface → ② configure
 87 the parallel environment number parameter (“env_num”) → ③ load the Verilog program onto the
 88 FPGA → ④ import the Python environments → ⑤ run the program. FPGA-Gym provides a set of reliable
 89 implementations of state-of-the-art RL algorithms. The FPGA_Gym library has included off-the-
 90 shelf RL environment examples. CPUs only handle data packing and transmission since the parallel
 91 environment computation runs entirely on the FPGA. This setup requires minimal dependencies, only
 92 using the “os” and “NumPy” [44] packages, making it highly compatible with various parallel RL
 93 algorithm libraries. Fig. 2(a) shows the typical source code to work with Stable-Baseline3 [45], a
 94 popular training framework for RL in Python, with the help of FPGA-Gym.

95 **Custom environment users**, who would like to customize new environments, need to ① implement
 96 the new environment step function in Verilog HDL [46] with the provided template (detailed in
 97 Section 3.3 and the supplementary materials) → ② update the basic environment computation module
 98 in the template and modify the parameters in the template (see the supplementary materials for
 99 step-by-step guidance).

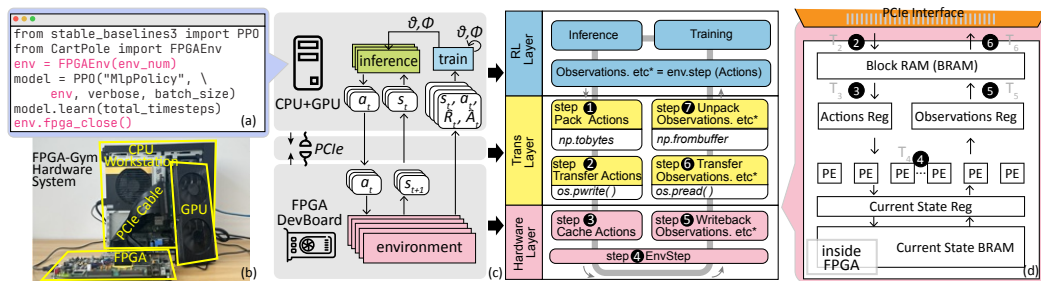


Figure 2: (a) Typical source code to run FPGA-Gym; (b) Photo of the proposed FPGA-Gym hardware system; (c) System architecture diagram of FPGA-Gym in a hierarchical manner. “Observations Reg” stores the observations, rewards, and terminate signals in FPGA computing rollouts. “PE” is the processing element synthesized in FPGA, which is a basic computing circuit module for a single environment instance.

100 3.2 Implementing FPGA-Gym framework

101 FPGA-Gym offloads the dominating environment step computation onto FPGA. The environment
 102 step in RL refers to the iterative process of obtaining the corresponding observation, reward, and
 103 termination based on the action generated by the agent. Each environment update is done in 7 steps
 104 in FPGA-Gym, as illustrated in Fig. 2(c): ①Pack data: CPU packs environment initial states and
 105 agent actions into bytes. ②Transfer actions: CPU transfers the packed data to FPGA block random

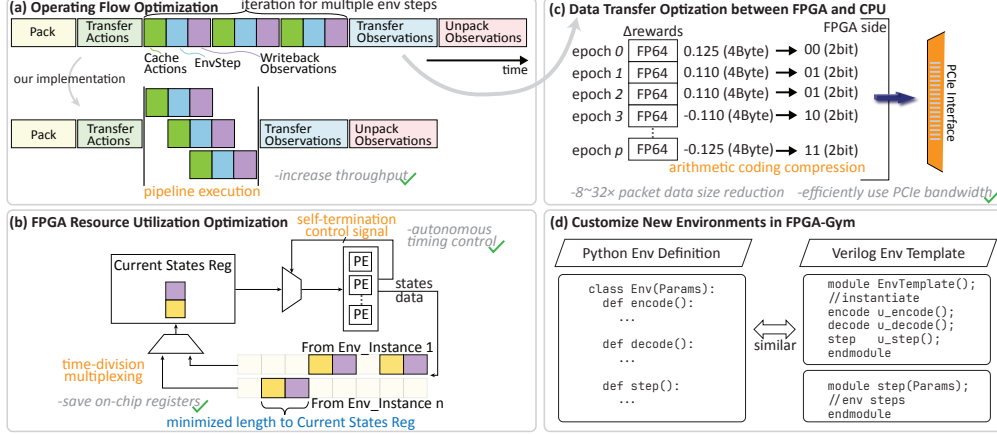


Figure 3: (a) The operating flow optimization; (b) FPGA resource utilization optimization; (c) data transfer optimization; (d) a template to customize new environments in FPGA-Gym. The orange text labels the proposed techniques.

106 access memory (BRAM) via PCIe. ③Cache actions: reading initial states, current states, actions, 107 and terminations from BRAM into registers for parallel computation. ④Environment step: when the 108 last step execution ends, the compute units use the initial states and actions as inputs. Otherwise, 109 they use the current states and actions as inputs. The compute units then output the next states, 110 observations, rewards, and termination flags. Then the current states are replaced by the next states. 111 ⑤Writeback observations: storing observations, rewards, and terminates into BRAM. ⑥Transfer 112 observations: transferring observations, rewards, and terminates from BRAM to CPU via PCIe. 113 ⑦Unpack: CPU unpacks observations, rewards, and terminates bytes into Python object/data types. 114 ①, ②, ⑥, and ⑦ is realized on the CPU side in Python. It seamlessly integrates with reinforcement 115 learning algorithms within Python code. ③, ④, ⑤ runs independently on FPGA and is implemented 116 with Verilog HDL [46]. To boost the computational throughput, we propose the following techniques 117 to achieve highly pipelined and compressed computation.

118 3.2.1 Pipeline design

119 To achieve acceleration, we first model and profile the aforementioned workflow latency step by step. 120 T_2 , T_3 , T_4 , T_5 and T_6 represent time costs of step ② to ⑥ depicted in Fig. 2(c), respectively:

$$\begin{aligned}
 T_2 &= \frac{k_1 \times env_num}{v_t} & T_3 &= \frac{k_1 \times env_num}{v_b} \\
 T_4^{without_pipeline} &= c \times env_num & T_4^{with_pipeline} &= c \\
 T_5 &= \frac{k_2 \times env_num}{v_b} & T_6 &= \frac{k_2 \times env_num}{v_t}
 \end{aligned} \tag{1}$$

121 where env_num is the number of environments to update, c (unit: second) is the time to compute 122 one environment step, k_1 (unit: GByte) is the data size of action per environment, k_2 (unit: GByte) is 123 the data size of observation, reward, and termination per environment, v_b (unit: GByte/second) is the 124 transfer speed between BRAM and register, v_t (unit: GByte/second) is transfer (PCIe) bandwidth 125 between CPU and FPGA. For steps ① and ⑦, we adopted efficient data packing and unpacking 126 strategies that consider both software and hardware aspects. Thanks to the robust functionality of 127 NumPy, these processes execute very quickly and occupy only a small fraction of the total runtime. 128 The measured and estimated time $T_{2\sim6}$ will be given in Section 4 and the supplementary materials.

129 We instantiate multiple environment compute units on FPGA so that they can execute the step of 130 multiple environments in a few clock cycles (often around 3~10 ns/cycle) at the same time. Instead 131 of realizing the aforementioned 7 steps sequentially, we implement a pipeline design specifically 132 between step ④ and step ⑤ so that the computing (④) and access to BRAM (⑤) inside FPGA occur 133 simultaneously to improve the computing throughput. Because the computation overlaps with the 134 data transfer in time, c (the time required to compute a single environment) and $(pe_num \cdot k_2/v_b)$

135 (the time required to transfer the computed data of a single environment to BRAM) should be equal
136 to each other to avoid the case that one step waiting for the other to complete.

137 3.2.2 Time-division multiplexing & self-termination reset

138 Storing a large number of current states would consume FPGA on-chip registers as the parallelism is
139 high. We carry out time-division reuse of this part of resources (Fig. 3(b)). Only when this part of
140 data is used, will it write in current states register from BRAM. This technique leads to an efficient
141 usage of FPGA on-chip computing and memory resources, i.e. the FPGA on-chip hardware resources
142 do not grow dramatically as the number of environments (“env_num”) increases. Additionally, in
143 a vectorized parallelism approach, managing individual environment resets upon the reception of
144 a “done” signal can be challenging. To solve this difficulty, we employ the self-termination reset in
145 the control flow. Each computation unit receives the action associated with a specific environment,
146 accompanied by the current state and termination signal from the preceding step, as well as the initial
147 state. If the termination signal is asserted, the environment automatically reverts to the initial state as
148 the input for the subsequent step calculation. This method voids significant wait times or processing
149 overhead.

150 3.2.3 Local access & data compression

151 The quantity of data transmitted has a direct impact on the values of the coefficients k_1 and k_2 in
152 Equation 1. To optimize data transmission, we employ the following techniques in FPGA-Gym. Firstly,
153 the state of the environment is preserved within the FPGA and updated automatically following each
154 calculation, except the initial state, which is received only once at the beginning. In subsequent
155 iterations, only the action sent from the CPU is accepted as an input. Secondly, for tasks where
156 the reward value encompasses several types but varies significantly, our framework uses arithmetic
157 coding to compress the reward values (illustrated in Fig. 3(c)). The reward values are replaced by
158 a categorical code, followed by batch numerical processing upon returning to the CPU. Lastly, for
159 tasks that allow for the concatenation of multiple values into a 32-bit number, we use bit-packing and
160 unpacking techniques to compress data, thereby enhancing data utilization efficiency. Details can be
161 found in the supplementary materials.

162 3.3 Customizing new environments

163 FPGA-Gym is a general RL acceleration framework that supports customized environments. We
164 have encapsulated the aforementioned workflow, pipeline, and compression into the backbone code
165 of FPGA-Gym. We provide a template for the customized environment step function (an epitome is
166 given in Fig. 3(d)). This source code template leaves all of the environment step functions empty and
167 the user can directly use Verilog HDL language (very similar to Python when it comes to the flow
168 control) to implement their environment or translate the environment Python code into Verilog HDL
169 code. Once done, the next step is to set up the key parameters in the template and synthesize digital
170 logic circuits on FPGA with the filled template.

171 4 Experiments

172 This section investigates the efficacy and compatibility of the FPGA-Gym framework for RL. In
173 particular, we would like to quantitatively answer the following questions: ① *Q1*: How effective
174 is FPGA-Gym in accelerating different types of RL environments? ② *Q2*: How significant are the
175 proposed design techniques (detailed in Section 3.2) on the parallel environment execution speedup?
176 ③ *Q3*: How to work with FPGA-Gym together the existing RL training algorithms and frameworks?

177 **Benchmarks & Baselines:** We select the typical environments *CartPole*, *Pendulum*, *CliffWalking*,
178 and *BlackJack* in Gymnasium [27] as the representatives for benchmarking. These 4 environments
179 cover discrete/continuous variable space, partially/globally observable simulation, and various cate-
180 gories including classic control, grid world, and strategy games (Table 1). The acceleration framework
181 baselines are EnvPool [22] and VectorEnv [27]. EnvPool is the latest and fastest (to the best of our

Table 1: Typical Environments For Benchmarking

Environment Name	Environment Type	Action Space Type	Observation Space Type	Rewards Type	Observability
<i>CartPole</i>	physical control	discrete	continuous	discrete	full
<i>Pendulum</i>	physical control	continuous	continuous	continuous	full
<i>MountainCar</i>	physical control	discrete	continuous	discrete	full
<i>CliffWalking</i>	gird world	discrete	discrete	discrete	full
<i>FrozenLake</i>	gird world	discrete	discrete	discrete	full
<i>Taxi</i>	gird world	discrete	discrete	discrete	full
<i>BlackJack</i>	strategy game	discrete	discrete	discrete	partial
<i>Seaquest</i>	atari game	discrete	discrete	discrete	full

182 knowledge) before this work to implement the above 4 environments. `VectorEnv` is a widely-used
 183 vectorized environment native to Gymnasium libraries using multithreading.

184 **Hardware Environment Setup:** The hardware platform we use has an off-the-shelf Xilinx VC707
 185 development board with 8-lane PCIe2.0 interfacing to a workstation with an Intel Core i9-10900K
 186 CPU and 128GB DDR4 memory. For the rollout length, 10^5 iterations are used to average the speed
 187 measurement experiments. All the training experiments set fixed 5 random seeds for reproducibility.
 188 For training, we employ DQN and PPO algorithms. PPO is implemented by directly calling the RL
 189 training programming framework Stable-Baseline3 [45, 47].

190 **4.1 Experiments on acceleration**

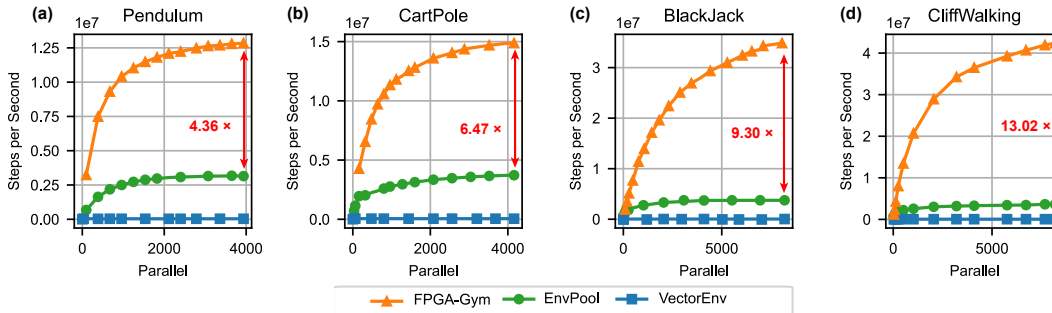


Figure 4: Measured throughputs of parallel environment steps (rollouts).

191 To answer *Q1*, we measure the throughputs at different parallelism of *Pendulum*, *CartPole*, *BlackJack*,
 192 and *CliffWalking* environments compared with `EnvPool` and `VectorEnv` frameworks, shown in
 193 Fig. 4. The x-axis is the number of instantiated parallel environments to compute. The y-axis is
 194 the total number of steps updated per second, which represents the computational throughput of
 195 executing environments in parallel (higher is better). Compared to `EnvPool`, `FPGA-Gym` improves
 196 the throughput by $4.36\times$ for *Pendulum*, $6.47\times$ for *CartPole* when executing 4000 instantiated
 197 environments in parallel, and $9.30\times$ for *BlackJack*, $13.02\times$ for *CliffWalking* with 8000 parallel
 198 environments. When running 2048 parallel *Seaquest* environments, `FPGA-Gym` achieves a $972.6\times$
 199 speedup compared to `EnvPool`. This validates the significant improvement of our approach to
 200 accelerate parallel RL environment execution. As the number of parallel environments increases,
 201 `EnvPool` reaches its maximum throughput earlier than `FPGA-Gym`, indicating a better scalability of
 202 `FPGA-Gym`. This confirms the original motivation of our proposed solution. The `EnvPool`'s curves in
 203 Fig. 4 reach saturation because of the limited number of computing threads; `FPGA-Gym`'s curves get
 204 lower slopes due to the BRAM bandwidth and PCIe transmission bandwidth.

205 **4.2 Experiments on pipeline-driven efficiency and hardware resource utilization**

206 To answer *Q2*, we test the execution latency breakdown before and after implementing pipeline
 207 optimization as well as the hardware resource utilization improvements. Fig. 5 shows the measured

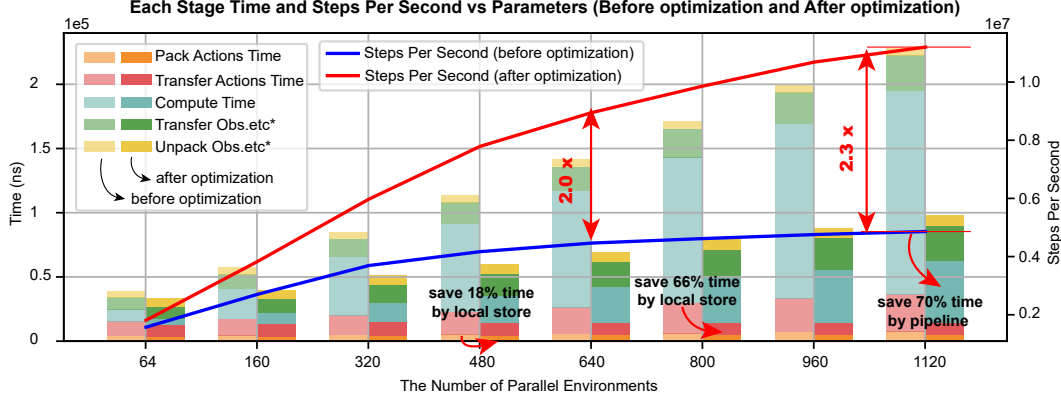


Figure 5: Each Stage Time and Steps Per Second vs The Number of Parallel Environments (Before and After Optimization) *: "Obs. etc" represents for the observations, rewards, and terminates.

208 computing latency of each stage in the workflow and steps per second before and after pipeline
 209 optimization. Here we take the *CartPole* task [27] as an example.

210 The x-axis is the number of parallel environments. The left y-axis is the execution time (stacked
 211 bar chart, lower values are better). The right y-axis is the steps per second (curves, higher is better).
 212 In each group of the stacked bars, the left (right) stacked bars are the execution time before (after)
 213 optimization. As the number of parallel environments increases, the computing module dominates
 214 the entire execution time. For the task with 1,120 instantiated *CartPole* environments in parallel, the
 215 pipeline technique saves up to 70% of compute time. Moreover, storing environment states locally
 216 saves CPU to FPGA’s transmission time by 66.0% and data packing time by 33.1%. With everything
 217 considered, FPGA-Gym achieves a $2.3\times$ speedup than before optimization.

218 It is convenient in FPGA-Gym to scale out the number of parallel environments. By tuning the
 219 “env_num” parameter, the FPGA-Gym framework can automatically complete the instantiation with
 220 more FPGA resources. Table 2 shows FPGA resource utilization at different parallel environment
 221 numbers. “pe_num” is the number of compute modules that compute the step of one environment.
 222 “c” is the time that a single RL environment computes. Look-up tables (LUTs), flip-flops (FFs)
 223 are the components to implement combinational logic, and registers to store intermediate results,
 224 respectively. The numbers in the last 3 rows are the FPGA implementation results for how many
 225 of the corresponding components are used. As shown, the LUTs and FFs of parallel *BlackJack*
 226 increase only 8.49% and 19.07% as env_num increases from 96 to 8000 ($84\times$), thanks to the proposed
 227 time-division multiplexing technique. Computing circuits (LUTs) and registers (FFs) can be shared
 228 among the environment instances.

Table 2: FPGA resource utilization

Environment	<i>Cliffwalking</i>		<i>Cartpole</i>		<i>Pendulum</i>		<i>BlackJack</i>	
env_num	64	8000	480	4000	192	4032	96	8064
pe_num / c (ns)	2/20		20/800		24/960		16/160	
power(W)	4.118	4.258	8.388	8.898	9.87	10.38	4.471	4.644
LUTs	27364	43481	175725	196896	228021	260307	46832	50811
FFs	26829	42821	254842	258832	267556	390253	41937	49938

229 4.3 Experiments on end-to-end agent training

230 To answer *Q3*, we evaluate the effects of FPGA-Gym on the existing RL algorithms and frameworks.
 231 We deploy the PPO algorithm from Stable-Baselines3 in 64 parallel *CartPole* environments with
 232 three frameworks—*VectorEnv*, *EnvPool*, and *FPGA-Gym*. The results are shown in Fig 6 (a) and
 233 (b). In Fig. 6, the x-axis of (a), (c), and (e) is the time of the entire RL training process, including

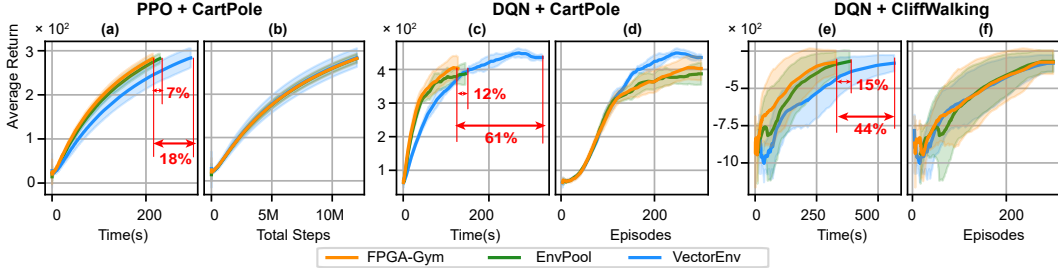


Figure 6: End_to_End_Train.

234 training, inference, and environment steps, etc. The x-axis of (b) is the total steps executed in the RL
 235 training process. The x-axis of (d), and (f) is the number of episodes executed during RL training.
 236 The y-axis is the measured average rewards during training. Fig. 6(a)(b) indicates that the training
 237 performance of the three methods is similar for the same number of steps. However, FPGA-Gym
 238 consistently achieves this performance faster than EnvPool and VectorEnv. When the total steps
 239 reach 12,000,000, the training time of FPGA-Gym is reduced by 18% and 7% compared to VectorEnv
 240 and EnvPool, respectively. Fig. 6(c)(d) presents the training performance of the *CartPole* using
 241 the DQN algorithm with 64 parallel environments. When episode number reaches 300, the training
 242 time of FPGA-Gym is reduced by 61% and 12% compared to VectorEnv and EnvPool, respectively.
 243 Fig. 6(e)(f) illustrates the training performance of the *Cliffwalking* environment using the DQN
 244 algorithm with 160 parallel environments. FPGA-Gym, EnvPool and VectorEnv achieve similar
 245 average rewards at the same episodes. When the number of episodes reaches 300, the training time of
 246 FPGA-Gym is reduced by 44% and 15% versus VectorEnv and EnvPool, respectively. The above
 247 results validate the compatibility of FPGA-Gym with various RL algorithmic libraries and frameworks,
 248 as well as the acceleration of training.

249 5 Limitation & Future Work

250 The primary motivation behind FPGA-Gym is to create a comprehensive hardware-enabled RL accel-
 251 eration framework that maximizes the potential of heterogeneous computing platforms, including
 252 CPUs, GPUs, and FPGAs. In its current form, FPGA-Gym supports a basic combination of GPU
 253 for neural networks, CPU for flow control, and FPGA for rollout acceleration. Our plans include
 254 exploring efficient task distribution and delicate scheduling among CPUs, GPUs, and FPGAs to
 255 enhance performance in dynamic, unpredictable open-world scenarios. Furthermore, we aim to
 256 enhance the scalability of FPGA-Gym by integrating a network of multiple FPGAs in upcoming
 257 versions. In the future, we plan to improve the BRAM bandwidth for better speedup ratios using
 258 multiple PCIe controllers and faster PCIe interfaces. With these enhancements, FPGA-Gym is set to
 259 deliver increased speed and scalability.

260 6 Conclusion

261 This study presents FPGA-Gym, a novel acceleration framework that combines FPGA and CPU
 262 components to enhance the performance of RL environments' parallel executions. By offloading
 263 environment steps to FPGA reconfigurable hardware, FPGA-Gym achieves a significant speedup of
 264 4.36 to 972.6 \times over the fastest existing software-based framework. We provide a modular design and
 265 parametric approach to simplify the deployment of environment updates for users without extensive
 266 FPGA expertise. A specialized FPGA-based parallel template is introduced for step operations in RL
 267 environments, similar to JAX's functionality but focusing on individual environments as the smallest
 268 units of computation. This approach offers flexibility and customization. The FPGA-Gym framework
 269 is compatible with existing RL algorithms and frameworks, including DQN and PPO algorithms. A
 270 standard environment library, similar to Gym library, is provided based on the FPGA-Gym framework.
 271 The library is publicly available https://github.com/Selinaee/FPGA_Gym.

272 Acknowledgments and Disclosure of Funding

273 This work was supported by National Natural Science Foundation of China (no. T2350006, 92264201,
274 92364102); and the 111 Project under Grant B18001. This work is sponsored by Beijing Nova
275 Program. The authors declare no competing interests.

276 References

- 277 [1] Zeyuan Ma, Hongshu Guo, Jiacheng Chen, Zhenrui Li, Guojun Peng, Yue-Jiao Gong, Yining
278 Ma, and Zhiguang Cao. MetaBox: A Benchmark Platform for Meta-Black-Box Optimization
279 with Reinforcement Learning. *Conference on Neural Information Processing Systems (NeurIPS)*,
280 2024.
- 281 [2] Zhecheng Yuan, Sizhe Yang, Pu Hua, Can Chang, Kaizhe Hu, and Huazhe Xu. RL-ViGen:
282 A Reinforcement Learning Benchmark for Visual Generalization. In *Conference on Neural*
283 *Information Processing Systems (NeurIPS)*, 2024.
- 284 [3] Lasse Espeholt, Raphaël Marinier, Piotr Stanczyk, Ke Wang, and Marcin Michalski. SEED
285 RL: Scalable and Efficient Deep-RL with Accelerated Central Inference. In *International*
286 *Conference on Learning Representations (ICLR)*, 2019.
- 287 [4] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang,
288 Melih Elibol, Zongheng Yang, William Paul, Michael I Jordan, et al. Ray: A distributed
289 framework for emerging {AI} applications. In *Symposium on Operating Systems Design and*
290 *Implementation (OSDI)*, 2018.
- 291 [5] Jiaming Ji, Borong Zhang, Jiayi Zhou, Xuehai Pan, Weidong Huang, Ruiyang Sun, Yiran Geng,
292 Yifan Zhong, Josef Dai, and Yaodong Yang. Safety gymnasium: A unified safe reinforcement
293 learning benchmark. *Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- 294 [6] Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam
295 Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl
296 with importance weighted actor-learner architectures. In *International conference on machine*
297 *learning (ICML)*, 2018.
- 298 [7] Xuehai Pan, Mickel Liu, Fangwei Zhong, Yaodong Yang, Song-Chun Zhu, and Yizhou Wang.
299 Mate: Benchmarking multi-agent reinforcement learning in distributed target coverage control.
300 In *Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- 301 [8] Steven Dalton and Iuri Frosio. Accelerating Reinforcement Learning through GPU Atari
302 Emulation. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- 303 [9] Misha Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang,
304 Lerrel Pinto, and Pieter Abbeel. URLB: Unsupervised Reinforcement Learning Benchmark. In
305 *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- 306 [10] Christopher Yeh, Victor Li, Rajeev Datta, Julio Arroyo, Nicolas Christianson, Chi Zhang, Yize
307 Chen, Mohammad Mehdi Hosseini, Azarang Golmohammadi, Yuanyuan Shi, et al. SustainGym:
308 Reinforcement Learning Environments for Sustainable Energy Systems. In *Conference on*
309 *Neural Information Processing Systems (NeurIPS)*, 2024.
- 310 [11] Denis Tarasov, Alexander Nikulin, Dmitry Akimov, Vladislav Kurenkov, and Sergey Kolesnikov.
311 Corl: Research-oriented deep offline reinforcement learning library. In *Conference on Neural*
312 *Information Processing Systems (NeurIPS)*, 2024.
- 313 [12] Yun Qu, Boyuan Wang, Jianzhun Shao, Yuhang Jiang, Chen Chen, Zhenbin Ye, Liu Linc, Yang
314 Feng, Lin Lai, Hongyang Qin, et al. Hokoff: Real Game Dataset from Honor of Kings and its
315 Offline Reinforcement Learning Benchmarks. In *Conference on Neural Information Processing*
316 *Systems (NeurIPS)*, 2024.

- 317 [13] Rong-Jun Qin, Xingyuan Zhang, Songyi Gao, Xiong-Hui Chen, Zewen Li, Weinan Zhang, and
318 Yang Yu. Neorl: A near real-world benchmark for offline reinforcement learning. In *Conference*
319 *on Neural Information Processing Systems (NeurIPS)*, 2022.
- 320 [14] Xiao-Yang Liu, Ziyi Xia, Jingyang Rui, Jiechao Gao, Hongyang Yang, Ming Zhu, Christina
321 Wang, Zhaoran Wang, and Jian Guo. FinRL-Meta: Market environments and benchmarks for
322 data-driven financial reinforcement learning. In *Conference on Neural Information Processing*
323 *Systems (NeurIPS)*, 2022.
- 324 [15] Hua Wei, Jingxiao Chen, Xiyang Ji, Hongyang Qin, Minwen Deng, Siqin Li, Liang Wang,
325 Weinan Zhang, Yong Yu, Liu Linc, et al. Honor of kings arena: an environment for generalization
326 in competitive reinforcement learning. In *Conference on Neural Information Processing Systems*
327 *(NeurIPS)*, 2022.
- 328 [16] Akhil Raj Baranwal, Salim Ullah, Siva Satyendra Sahoo, and Akash Kumar. *ReLAccS* : A
329 multilevel approach to accelerator design for reinforcement learning on FPGA-based systems.
330 *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (IEEE TCAD)*,
331 2021.
- 332 [17] Yuan Meng, Sanmukh Kuppannagari, and Viktor Prasanna. Accelerating Proximal Policy
333 Optimization on CPU-FPGA Heterogeneous Platforms. In *International Symposium on Field-*
334 *Programmable Custom Computing Machines (FCCM)*, 2020.
- 335 [18] Hyungmin Cho, Pyeongseok Oh, Jiyoung Park, Wookeun Jung, and Jaejin Lee. FA3C: FPGA-
336 Accelerated Deep Reinforcement Learning. In *Conference on Architectural Support for Pro-*
337 *gramming Languages and Operating Systems (ASPLOS)*, 2019.
- 338 [19] Chan-Wei Hu, Jiang Hu, and Sunil P. Khatri. TD3lite: FPGA acceleration of reinforcement
339 learning with structural and representation optimizations. In *Conference on Field-Programmable*
340 *Logic and Applications (FPL)*, 2022.
- 341 [20] Yuan Meng, Chi Zhang, and Viktor Prasanna. FPGA acceleration of deep reinforcement learning
342 using on-chip replay management. In *Proceedings of the 19th ACM International Conference*
343 *on Computing Frontiers (CF)*, 2022.
- 344 [21] Samuel Wiggins, Yuan Meng, Rajgopal Kannan, and Viktor Prasanna. Accelerating multi-agent
345 DDPG on CPU-FPGA heterogeneous platform. In *IEEE High Performance Extreme Computing*
346 *Conference (HPEC)*, 2023.
- 347 [22] Jiayi Weng, Min Lin, Shengyi Huang, Bo Liu, Denys Makoviichuk, Viktor Makoviyichuk,
348 Zichen Liu, Yufan Song, Ting Luo, Yukun Jiang, Zhongwen Xu, and Shuicheng Yan. Envpool:
349 A Highly Parallel Reinforcement Learning Environment Execution Engine. In *Conference on*
350 *Neural Information Processing Systems (NeurIPS)*, 2022.
- 351 [23] Todd Hester, Matej Vecerík, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Horgan,
352 John Quan, Andrew Sendonaris, Ian Osband, Gabriel Dulac-Arnold, John P. Agapiou, Joel Z.
353 Leibo, and Audrunas Gruslys. Deep Q-learning From Demonstrations. In *AAAI Conference on*
354 *Artificial Intelligence (AAAI)*, 2018.
- 355 [24] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
356 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
357 Human-level control through deep reinforcement learning. *Nature*, 2015.
- 358 [25] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
359 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 360 [26] Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu.
361 The surprising effectiveness of ppo in cooperative multi-agent games. In *Conference on Neural*
362 *Information Processing Systems (NeurIPS)*, 2022.

- 363 [27] Mark Towers, Jordan K. Terry, Ariel Kwiatkowski, John U. Balis, Gianluca de Cola, Tristan
364 Deleu, Manuel Goulão, Andreas Kallinteris, Arjun KG, Markus Krimmel, Rodrigo Perez-
365 Vicente, Andrea Pierré, Sander Schulhoff, Jun Jet Tai, Andrew Tan Jin Shen, and Omar G.
366 Younis. Gymnasium. <https://zenodo.org/record/8127025>, 2023.
- 367 [28] Aleksei Petrenko, Zhehui Huang, Tushar Kumar, Gaurav S. Sukhatme, and Vladlen Koltun.
368 Sample Factory: Egocentric 3D Control from Pixels at 100000 FPS with asynchronous rein-
369 forcement learning. In *International conference on machine learning (ICML)*, 2020.
- 370 [29] Robert Tjarko Lange. gymmax: A JAX-based reinforcement learning environment library.
371 <http://github.com/RobertTLange/gymmax>, 2022.
- 372 [30] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based
373 control. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2012.
- 374 [31] Matteo Hessel, Manuel Kroiss, Aidan Clark, Iurii Kemaev, John Quan, Thomas Keck, Fabio
375 Viola, and Hado van Hasselt. Podracer architectures for scalable reinforcement learning. *arXiv*
376 *preprint arXiv:2104.06272*, 2021.
- 377 [32] Tian Lan, Sunil Srinivasa, Huan Wang, and Stephan Zheng. Warpdrive: Fast end-to-end deep
378 multi-agent reinforcement learning on a GPU. *Journal of Machine Learning Research (JMLR)*,
379 2022.
- 380 [33] Michael Matthews, Michael Beukman, Benjamin Ellis, Mikayel Samvelyan, Matthew Jackson,
381 Samuel Coward, and Jakob Foerster. Craftax: A lightning-fast benchmark for open-ended
382 reinforcement learning. *arXiv preprint arXiv:2402.16801*, 2024.
- 383 [34] Sotetsu Koyamada, Shinri Okano, Soichiro Nishimori, Yu Murata, Keigo Habara, Haruka
384 Kita, and Shin Ishii. Pgx: Hardware-Accelerated Parallel Game Simulators for Reinforcement
385 Learning. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- 386 [35] Clément Bonnet, Daniel Luo, Donal Byrne, Shikha Surana, Vincent Coyette, Paul Duckworth,
387 Laurence I. Midgley, Tristan Kalloniatis, Sasha Abramowitz, Cemlyn N. Waters, Andries P. Smit,
388 Nathan Grinsztajn, Ulrich A. Mbou Sob, Omayma Mahjoub, Elshadai Tegegn, Mohamed A.
389 Mimouni, Raphaël Boige, Ruan de Kock, Daniel Furelos-Blanco, Victor Le, Arnu Pretorius, and
390 Alexandre Laterre. Jumanji: a diverse suite of scalable reinforcement learning environments in
391 JAX. In *International Conference on Learning Representations (ICLR)*, 2023.
- 392 [36] Alexander Rutherford, Benjamin Ellis, Matteo Gallici, Jonathan Cook, Andrei Lupu, Garðar
393 Ingvarsson, Timon Willi, Akbir Khan, Christian Schröder de Witt, Alexandra Souly, Saptarashmi
394 Bandyopadhyay, Mikayel Samvelyan, Minqi Jiang, Robert Tjarko Lange, Shimon Whiteson,
395 Bruno Lacerda, Nick Hawes, Tim Rocktäschel, Chris Lu, and Jakob N. Foerster. Jaxmarl: Multi-
396 agent RL environments and algorithms in JAX. In *Proceedings of International Conference on*
397 *Autonomous Agents and Multiagent Systems (AAMAS)*, 2024.
- 398 [37] Mathias Lechner, Lianhao Yin, Tim Seyde, Tsun-Hsuan Johnson Wang, Wei Xiao, Ramin M.
399 Hasani, Joshua Rountree, and Daniela Rus. Gigastep - One Billion Steps per Second Multi-agent
400 Reinforcement Learning. In *Conference on Neural Information Processing Systems (NeurIPS)*,
401 2023.
- 402 [38] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles
403 Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac
404 Gym: High Performance GPU Based Physics Simulation For Robot Learning. In *Conference*
405 *on Neural Information Processing Systems (NeurIPS)*, 2021.
- 406 [39] C. Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier
407 Bachem. Brax - A differentiable physics engine for large scale rigid body simulation. In
408 *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.

- 409 [40] Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, Viacheslav Sinii, Artem Agarkov, and
410 Sergey Kolesnikov. XLand-minigrid: Scalable meta-reinforcement learning environments in
411 JAX. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- 412 [41] Benjamin Ellis, Jonathan Cook, Skander Moalla, Mikayel Samvelyan, Mingfei Sun, Anuj Ma-
413 hajan, Jakob Foerster, and Shimon Whiteson. Smacv2: An improved benchmark for cooperative
414 multi-agent reinforcement learning. In *Conference on Neural Information Processing Systems*
415 *(NeurIPS)*, 2024.
- 416 [42] Yiheng Zhu, Yang Zhan, Xuankun Huang, Yuwei Chen, Jiangwen Wei, Wei Feng, Yinzhi Zhou,
417 Haoyuan Hu, Jieping Ye, et al. Ofcourse: A multi-agent reinforcement learning environment for
418 order fulfillment. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2024.
- 419 [43] Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo Perez-Vicente, Lucas Willems,
420 Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld:
421 Modular & customizable reinforcement learning environments for goal-oriented tasks. In
422 *Conference on Neural Information Processing Systems (NeurIPS)*, 2024.
- 423 [44] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen,
424 David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern,
425 Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime
426 Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard,
427 Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant.
428 Array programming with NumPy. *Nature*, 2020.
- 429 [45] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
430 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*
431 *Machine Learning Research (JMLR)*, 2021.
- 432 [46] IEEE. IEEE Standard for Verilog Hardware Description Language. *IEEE Std 1364-2005*
433 *(Revision of IEEE Std 1364-2001)*, 2006.
- 434 [47] Antonin Raffin. RL Baselines3 Zoo. <https://github.com/DLR-RM/rl-baselines3-zoo>,
435 2020.