OCTAX: ACCELERATED CHIP-8 ARCADE ENVIRON-MENTS FOR REINFORCEMENT LEARNING IN JAX

Anonymous authorsPaper under double-blind review

ABSTRACT

Reinforcement learning (RL) research requires diverse, challenging environments that are both tractable and scalable. While modern video games may offer rich dynamics, they are computationally expensive and poorly suited for large-scale experimentation due to their CPU-bound execution. We introduce OCTAX, a high-performance suite of classic arcade game environments implemented in JAX, based on CHIP-8 emulation, a predecessor to Atari, which is widely adopted as a benchmark in RL research. OCTAX provides the JAX community with a long-awaited end-to-end GPU alternative to the Atari benchmark, offering imagebased environments, spanning puzzle, action, and strategy genres, all executable at massive scale on modern GPUs. Our JAX-based implementation achieves ordersof-magnitude speedups over traditional CPU emulators while maintaining perfect fidelity to the original game mechanics. We demonstrate OCTAX's capabilities by training RL agents across multiple games, showing significant improvements in training speed and scalability compared to existing solutions. The environment's modular design enables researchers to easily extend the suite with new games or generate novel environments using large language models, making it an ideal platform for large-scale RL experimentation.

1 Introduction

Modern reinforcement learning (RL) research (Sutton & Barto, 2018) demands extensive experimentation to achieve statistical validity, yet computational constraints severely limit experimental scale. RL papers routinely report results with fewer than five random seeds due to prohibitive training costs (Henderson et al., 2018; Colas et al., 2018; Agarwal et al., 2021; Mathieu et al., 2023; Gardner et al., 2025). While understandable from a practical standpoint, this undersampling undermines statistical reliability and impedes algorithmic progress. Environment execution creates this bottleneck: while deep learning has embraced end-to-end GPU acceleration, RL environments remain predominantly CPU-bound. Originally designed under severe hardware constraints, classic arcade games represent a solution for scalable RL experimentation. The Atari Learning Environment (ALE) (Bellemare et al., 2013) has established itself as a standard RL benchmark, although existing implementations remain fundamentally CPU-bound. As noted by Obando-Ceron & Castro (2020), the Rainbow paper (Hessel et al., 2018) required 34,200 GPU hours (equivalent to 1,425 days) of experiments, a computational cost that is prohibitively high for small research laboratories. In this paper, we propose an alternative approach for training RL agents in environments with mechanisms similar to ALE, with significantly reduced computational cost.

Contributions. We introduce OCTAX¹, a suite of arcade game environments implemented in JAX (Bradbury et al., 2018a) through CHIP-8 emulation. CHIP-8, a 1970s virtual machine specification contemporary with early Atari systems, became the foundation for numerous classic games spanning puzzle, action, and strategy genres. CHIP-8's constraint-driven design creates games with similar cognitive demands to Atari while enabling efficient vectorized emulation that scales to thousands of parallel instances. The JAX ecosystem has rapidly emerged as a solution for scalability in RL research but lacks native environments, particularly image-based ones. Our framework addresses

 $^{^1\}mbox{The anonymized repository containing all source code, experiments, and data is available at: https://anonymous.4open.science/r/octax-C8E8/README.md$

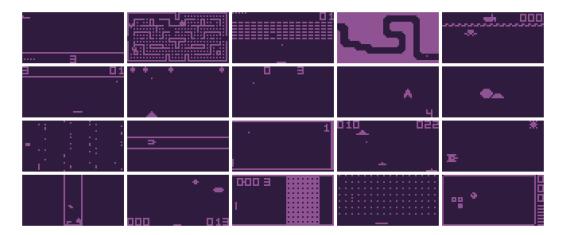


Figure 1: Overview of CHIP-8 game environments implemented in OCTAX.

this gap by transforming classic games into fully vectorized, GPU-accelerated simulations. These simulations run thousands of game instances in parallel while maintaining perfect fidelity to the original mechanics. This approach dramatically reduces experiment times. Experiments that previously required days or weeks can now be completed in hours. This efficiency makes comprehensive hyperparameter sweeps and ablation studies computationally feasible. The modular design facilitates extension with new games or automated generation using large language models that can directly output CHIP-8 assembly code. Figure 1 provides an overview of the integrated CHIP-8 games.

Outline. First, we present the our end-to-end JAX implementation of classic arcade environments through CHIP-8 emulation (Section 3). Second, we demonstrate diverse learning dynamics through PPO evaluation across 16 games (Section 4.1). Third, we achieve 350,000 environment steps per second (1.4 million frames per second) on consumer-grade hardware, substantially outperforming CPU-based solutions (Section 4.2). Fourth, we establish an LLM-assisted pipeline for automated environment generation that creates meaningful difficulty gradients (Section 4.3).

2 RELATED WORK

Game environments have proven essential for RL research because they provide engaging, humanrelevant challenges with clear success metrics. The Arcade Learning Environment (ALE) Bellemare et al. (2013) demonstrated this principle by establishing Atari 2600 games as the standard RL benchmark, enabling breakthrough algorithms like DQN (Mnih et al., 2015) and Rainbow (Hessel et al., 2018). The success of these classic arcade games stems from their constraint-driven design: simple rules that yield complex behaviors, deterministic dynamics that enable reproducible experiments, and visual complexity that tests spatial reasoning without overwhelming computational resources.

While algorithmic advances demand increasingly large-scale experiments with thousands of parallel environments and extensive hyperparameter sweeps, traditional game environments remain CPU-bound and poorly suited for parallel execution. This mismatch has driven a progression of solutions, each addressing different aspects of the scalability problem.

Game-based RL environment platforms. Increasingly sophisticated gaming platforms have been developed to test different dimensions of learning performance. NetHack Learning Environment (Küttler et al., 2020) provides procedurally generated roguelike challenges that test long-term planning, while Crafter (Hafner, 2021) offers simplified Minecraft-like environments focused on resource management. These environments expand cognitive challenges beyond arcade games, but their CPU-based implementations compound the scalability problem.

CPU high-performance solutions. Several projects have focused on optimizing CPU-based environment execution. EnvPool (Weng et al., 2022) achieves substantial speed improvements through highly optimized C++ implementation, demonstrating up to 1 million Atari frames per second on

high-end hardware. PufferLib (Suarez, 2025) provides environments written entirely in C, achieving millions of steps per second through over 20,000 lines of optimized code. While these approaches improve CPU throughput, they retain fundamental limitations: costly CPU-GPU data transfers during training and require C implementation in a Python-dominated field.

GPU-accelerated RL environments. GPU-accelerated solutions target the constraint more directly by moving environment execution to accelerators. CUDA Learning Environment (CuLE) (Dalton et al., 2020) provides a pioneering CUDA port of ALE, achieving 40-190 million frames per hour on single GPUs. Isaac Gym (Makoviychuk et al., 2021) demonstrates similar principles for robotics tasks, achieving 2-3 orders of magnitude speedups over CPU approaches by running thousands of environments simultaneously. These GPU approaches solve computational bottlenecks but introduce NVIDIA hardware dependence and substantial per-environment engineering costs.

JAX-based environments. The adoption of JAX (Bradbury et al., 2018b) has enabled natively accelerated environments that combine portability across hardware with end-to-end GPU acceleration. Brax (Freeman et al., 2021) established viability through MuJoCo-like physics simulation, while Gymnax (Lange, 2022) provides JAX implementations of classic control tasks and simplified environments from BSuite (Osband et al., 2019) and MinAtar (Young & Tian, 2019). Specialized environments target specific research needs: XLand-MiniGrid (Nikulin et al., 2024) and Navix (Pignatelli et al., 2024) focus on gridworld navigation, Jumanji (Bonnet et al., 2023) spans domains from simple games to NP-hard combinatorial problems, Pgx (Koyamada et al., 2023) provides classic board games, and PuzzleJAX Earle et al. (2025) enables dynamic compilation of puzzle games.

Despite this coverage, a critical gap remains: classic arcade games. While MinAtar provides simplified versions of Atari games, the full visual complexity and authentic game mechanics of classic arcade games remain absent from the JAX ecosystem. OCTAX addresses this gap by providing the first end-to-end JAX implementation of classic arcade games through CHIP-8 emulation, delivering computational benefits while preserving the engaging gameplay mechanics that made arcade games valuable for algorithmic development.

3 OCTAX: THE ACCELERATED CHIP-8 PLATFORM

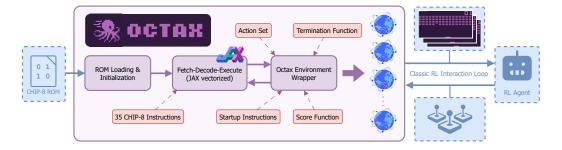


Figure 2: OCTAX architecture: ROM loading, CHIP-8 emulation pipeline, and RL environment integration. The system transforms game ROMs through fetch-decode-execute cycles into vectorized JAX operations suitable for GPU acceleration.

This section presents our JAX implementation of CHIP-8 emulation. We detail the design decisions that enable GPU acceleration while maintaining behavioral fidelity to original games, and explain how CHIP-8's architecture provides an optimal foundation for scalable experimentation in RL. Figure 2 summarizes this section.

3.1 Why CHIP-8 for RL research?

CHIP-8 represents a strategic choice for RL environment design. Created in the 1970s as a virtual machine specification, CHIP-8 features a 64×32 monochrome display, 16 registers, 4KB memory, and 35-instruction set. These constraints, originally imposed by early microcomputer limitations, create several research advantages.

The platform provides image-based environments comparable to Atari games while offering some computational advantages. The 4KB memory footprint allows thousands of simultaneous game instances without memory constraints. The simple instruction set reduces emulation overhead compared to complex modern processors. The deterministic execution model ensures experimental reproducibility across different hardware configurations.

The platform supports everything from precise action games requiring split-second timing to complex puzzles demanding long-horizon planning. The 16-key input system provides sufficient complexity for interesting control challenges while remaining tractable for systematic analysis. Most importantly, CHIP-8 games are inherently modifiable and analyzable: their simple assembly code can be automatically generated, modified, and assessed for difficulty, enabling novel research directions in environment design and curriculum learning. This combination of Atari-like visual complexity with modern computational efficiency makes CHIP-8 well-suited for the JAX ecosystem, where extensive parallelization can transform week-long experiments into hour-long runs.

3.2 How does Octax work?

OCTAX converts CHIP-8 ROMs² into vectorized RL environments while maintaining compatibility with original games. The implementation leverages JAX's functional programming model and vectorization capabilities to enable GPU acceleration.

ROM loading and initialization. Game data is loaded from .ch8 files into the emulator's 4KB memory space starting at address 0x200, following the standard CHIP-8 program layout first introduced in Weisbecker (1978). The system initializes with font data at address 0x50, sixteen general-purpose registers (V0-VF), an index register (I), a program counter (PC), and the 64×32 monochrome display buffer.

Fetch-decode-execute cycle. The core emulation loop implements the classic processor cycle using JAX primitives. The $\mathtt{fetch}()$ function retrieves 16-bit instructions from memory and advances the program counter. The $\mathtt{decode}()$ function extracts instruction components through bitwise operations, identifying opcodes, register indices, and immediate values. The $\mathtt{execute}()$ function uses JAX's $\mathtt{lax.switch}$ for GPU-compatible instruction dispatch to specialized handlers.

Vectorized instruction execution. Instruction handlers follow JAX's functional programming model, treating state as immutable and returning updated copies. ALU operations handle arithmetic and bitwise logic with carry/borrow flag management. Control flow instructions implement jumps, calls, and conditional operations using lax.cond. The display system uses vectorized operations to render sprites across the entire framebuffer simultaneously.

Environment integration. The OctaxEnv wrapper transforms the emulator into a standard RL interface. Each RL step executes multiple CHIP-8 instructions to maintain authentic game timing relative to the original 700Hz instruction frequency. The default frame skip setting preserves realistic game dynamics. Observations consist of the 64×32 display with 4-frame stacking, producing (4, 64, 32) boolean arrays. Actions map from discrete RL outputs to game-specific key subsets plus a no-op option. The wrapper manages delay and sound timers at 60Hz and executes startup sequences to bypass menu screens.

3.3 How to transform games into RL environments?

Converting CHIP-8 games into RL environments requires extracting reward signals and termination conditions from game-specific memory layouts and register usage patterns.

Score function design. Games store scores in different registers using various encoding schemes. OCTAX provides game-specific <code>score_fn</code> functions that extract scores from appropriate memory locations. Brix stores its score in register V5, incrementing with each destroyed brick. Pong encodes scores in BCD format within register V14, requiring <code>score = (V[14] // 10) - (V[14] % 10)</code> to compute player advantage. This flexibility allows researchers to experiment with alternative reward formulations based on different game state components.

²ROM stands for Read-Only Memory, a type of storage originally used in game cartridges to hold software that cannot be modified by the user.

Termination logic. Games signal completion through different register states that must be identified through analysis. Brix terminates when lives (V14) reach zero, while Tetris uses a dedicated game-state register (V1) that equals 2 on game over. Some games require compound conditions: Space Flight ends when either lives reach zero or a level completion counter exceeds a threshold, implemented as terminated = $(V[9] = 0) | (V[12] >= 0 \times 3E)$.

Action space optimization. Most games use subsets of the 16-key hexadecimal keypad. OCTAX supports custom action_set arrays that map RL action indices to relevant keys. Pong requires only keys 1 and 4 for paddle movement, while Worm uses directional keys 2, 4, 6, 8. This reduces action space size and accelerates learning by eliminating irrelevant inputs.

Initialization handling. Many games include menu screens that interfere with RL training. OCTAX supports startup_instructions parameters that automatically execute instruction sequences during environment reset, bypassing menus to begin gameplay immediately.

We address CHIP-8's non-standardized scoring and termination by combining static ROM analysis and dynamic memory monitoring during gameplay, as detailed in Appendix C.

3.4 WHICH GAMES DOES OCTAX SUPPORT?

OCTAX provides a curated collection of classic CHIP-8 games across multiple genres and difficulty levels. The current implementation includes 21 titles, with additional games planned for future releases. All environments maintain full compatibility with both Gymnasium and Gymnax APIs.

Category	Available Games	Required Capabilities
Puzzle	Tetris, Blinky, Worm	Long-horizon planning, spatial reasoning
Action	Brix, Pong, Squash, Vertical Brix, Wipe Off, Filter	Timing, prediction, reactive control
Strategy	Missile Command, Rocket, Submarine, Tank Battle, UFO	Resource management, tactical decisions
Exploration	Cavern (7 levels), Flight Runner, Space Flight (10 levels), Spacejam!	Spatial exploration, continuous navigation
Shooter	Airplane, Deep8, Shooting Stars	Simple reaction, basic timing

Table 1: Currently implemented games in OCTAX.

The games (Figure 1) vary across multiple dimensions of difficulty and cognitive demand. Temporal complexity ranges from immediate reactions to long-term planning requirements. Spatial complexity spans single-screen environments to multi-screen worlds requiring navigation. Reward structures include both dense scoring mechanisms and sparse achievement-based systems. This systematic variation enables controlled studies of algorithmic performance across different challenge types while maintaining a unified technical framework for fair comparison. A categorization of these games is provided in Table 1, with more detailed descriptions available in Appendix C.2.

4 EXPERIMENTAL EVALUATION

We evaluate OCTAX through RL training experiments across 16 diverse CHIP-8 games. Our goal is to demonstrate that the environments exhibit varied difficulties and learning dynamics suitable for RL research and benchmark the platform's computational performance.

4.1 How do RL agents learn in Octax?

We train Proximal Policy Optimization (PPO) (Schulman et al., 2017) agents across our game suite due to its widespread adoption and proven scalability with parallel environments (Rudin et al., 2022).

Network architecture. Our PPO agent³ uses a convolutional neural network designed for OCTAX's (4, 64, 32) stacked observations. The feature extractor consists of three convolutional layers

³Based on Rejax implementation (Liesen et al., 2024).

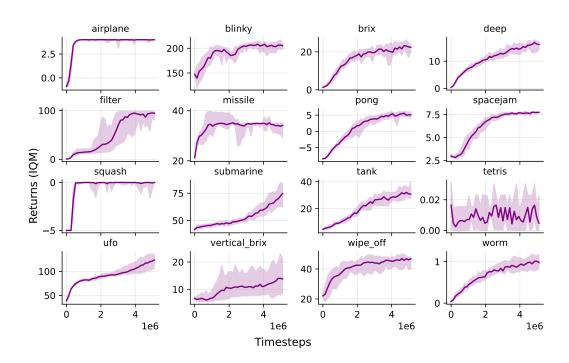


Figure 3: PPO learning curves across 16 games: Interquartile Mean (IQM) returns using 10th-90th percentile ranges over 5M timesteps, with confidence intervals computed across 12 random seeds.

with 32, 64, and 64 filters respectively. These layers use kernel sizes of (8,4), 4, and 3 with corresponding strides of (4,2), 2, and 1. Extracted features are flattened and fed to separate actor and critic heads, each containing a single 256-unit hidden layer with ReLU activation throughout.

Training configuration. We combine grid search optimization (detailed in Appendix 4) on Pong with CleanRL's standard Atari PPO hyperparameters (Huang et al., 2022). This yields GAE lambda of 0.95, clipping epsilon of 0.2, value function coefficient of 0.5, and entropy coefficient of 0.01. Each experiment uses 512 parallel environments with 32-step rollouts, 4 training epochs per update, and 32 minibatches for gradient computation. We apply the Adam optimizer (Kingma & Ba, 2014) with learning rate 5×10^{-4} and gradient clipping for stable training across 5 million timesteps per environment.

Experimental setup. We conduct 12 independent training runs per game using different random seeds. All experiments run on a single NVIDIA A100 GPU with 24 concurrent training sessions. Agent performance is assessed every 131,072 timesteps on 128 parallel environments.

Results analysis. The training curves in Figure 3 reveal distinct learning profiles across games. We observe three main patterns that reflect different cognitive demands. *Rapid plateau games* (Airplane, Brix, Deep, Filter, Blinky) show quick initial learning followed by stable performance, suggesting clear reward signals. *Gradual improvement games* (Submarine, Tank, UFO) exhibit sustained learning throughout training, indicating either sparser reward structures or more complex strategic requirements. *Limited performance games*, like Tetris, exhibit significant variance with little absolute progress, making them difficult for standard policy gradient methods. Similarly, in Worm (a Snake clone), agents often manage to eat only a single apple before dying.

These learning profiles support the cognitive diversity of CHIP-8 environments, demonstrating that different games test varied aspects of learning and control. Individual training runs averaged 65 minutes each, with 24 experiments running concurrently, achieving approximately 30,800 environment steps per second across all parallel sessions.

4.2 How does Octax scale with parallelization?

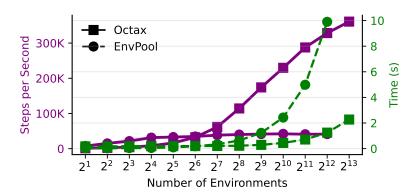


Figure 4: Performance scaling of OCTAX and EnvPool across parallelization levels. The solid purple line is the number of steps per second (higher is better), and the dashed green line is the total execution time in seconds (lower is better).

Experimental setup. We measure environment throughput across different parallelization levels to quantify OCTAX's computational advantages. This experiment isolates pure computational benefits by fixing the game (Pong) and agent behavior (constant action) while varying parallel environment instances. Since all environments execute identical CHIP-8 computational cycles, these performance measurements apply uniformly across the entire game suite. To better interpret our results, we compare against EnvPool because it is widely adopted in RL research, using ALE Pong to assess CPU vs. GPU-based environment scalability.

Configuration. We benchmark on a consumer-grade workstation with an RTX 3090 (24GB VRAM), 32GB RAM, and an Intel i7 processor (20 cores). We measure execution time for 100-step rollouts across varying parallel environment counts, with 50 independent measurements per configuration. The primary metric is environment steps per second, calculated as (number of environments × 100 steps) divided by execution time, where each step represents 4 frames due to OCTAX's default frame skip setting.

Performance results. Figure 4 demonstrates near-linear scaling up to 350,000 steps (or 1,4M frames) per second with 8,192 parallel environments before hitting VRAM limitations. EnvPool running ALE Pong with all available CPU cores shows reduced scaling, plateauing around 25,000 steps per second due to CPU saturation. OCTAX achieves a 14× improvement in computational efficiency at high parallelization levels, reducing the computational cost of large-scale RL experiments. We also measured GPU memory usage across different environment counts, finding that execution memory scales linearly with the number of parallel environments with our benchmark script, consuming approximately 2 MB of GPU memory per environment

4.3 How do LLMs assist environment creation?

Large language models (LLMs) have demonstrated a strong capability in code generation across diverse programming languages, enabling the automated creation of environments in RL research. Here we explore OCTAX's capacity to accelerate research by leveraging LLMs to generate novel tasks, extending beyond manually designed game suites toward automated environment synthesis, as explored in Faldor et al. (2024).

Context. During OCTAX's development, we encountered several games where reward and termination logic proved difficult to extract through manual analysis of game mechanics. In these cases, we decompiled ROMs to obtain CHIP-8 assembly code and successfully employed LLMs to generate appropriate score_fn and terminated_fn functions by analyzing the assembly instructions. This process revealed LLMs' capability to understand low-level game logic and translate it into RL-compatible reward structures. This success motivated us to investigate the reverse

pipeline: using LLMs to generate complete CHIP-8 games from high-level descriptions, then leveraging OCTAX's scalable simulation to evaluate these procedurally created environments.

Automated environment generation pipeline. Our pipeline consists of seven replicable steps for automatic CHIP-8 game generation. In Step 1, we construct a corpus of CHIP-8 tutorials, documentation, and programming examples, ensuring the LLM understands the architecture's instruction set, memory layout, and common coding patterns. In Step 2, we embed this corpus into a prompt (detailed in Appendix E.1) that guides the LLM to produce syntactically correct CHIP-8 programs from high-level instructions. In Step 3, we provide a description of the game with desired mechanics, objectives, and constraints. In Step 4, the LLM generates the initial CHIP-8 code based on the provided description. In Step 5, an automated feedback loop between the LLM and a CHIP-8 compiler iteratively refines the code based on compilation errors until successful. In Step 6, Python wrapper functions for score_fn and terminated_fn are automatically generated, translating CHIP-8 registers into RL-compatible reward and termination signals. Finally, in Step 7, the game description is augmented to increase difficulty or introduce new challenges. Both the new description and the previously generated game are added to the LLM's context before next iteration. Figure 5 summarizes the automated environment generation pipeline.

Target Shooter case study. We validated this pipeline using Claude Opus 4.1, known for its proficiency in programming, with the following description: "Target Shooter – Targets appear randomly on the screen, and the player moves a crosshair to shoot them. Score increases per hit, and the game ends after a fixed number of targets." The system successfully generated three progressive difficulty levels: static targets for basic aiming skills, time-limited targets introducing decision pressure, and moving targets with time constraints requiring predictive aiming. Each level maintains consistent register mappings for score and termination, simplifying OCTAX compatibility. Figure 6 shows how the LLM-generated environment visual appearance. All the code generated by the LLM is given in E.2,

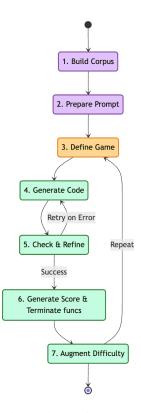


Figure 5: Environment generation pipeline.

RL experiments. Using identical PPO configurations from Section 4.1, we trained agents on the three generated difficulty levels over 5M timesteps. Figure 7 demonstrates clear performance stratification across difficulty levels: Level 1 agents achieved optimal returns of 10.0 with rapid convergence by 1M timesteps, Level 2 agents plateaued at 9.0 returns with moderate learning speed, while Level 3 agents reached 8.0 returns with the slowest progression. The inverse relationship between difficulty level and both final performance and sample efficiency indicates that our LLM-generated environments successfully create a meaningful difficulty gradient. This proof-of-concept demonstrates the feasibility of automated environment generation for RL research via OCTAX, with promising applications in curriculum learning, open-endedness, and continual learning scenarios.



Figure 6: Rendering of the Target Shooter game showing the player (left, circular object) and target (right, cross-shaped object).

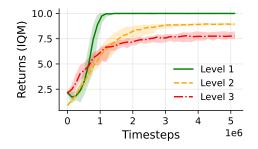


Figure 7: PPO training performance on generated environments with varying difficulty.

5 CONCLUSION

We introduced OCTAX, a JAX-based CHIP-8 emulation platform that provides GPU-accelerated arcade game environments for reinforcement learning research. Our implementation achieves significant performance improvements over CPU-based alternatives, enabling experiments with thousands of parallel environments while maintaining perfect behavioral fidelity to original games. Through PPO evaluation across 16 diverse games, we demonstrated varied learning dynamics that highlight the cognitive diversity within classic arcade environments. The platform's modular design enables both manual game integration and automated environment generation using large language models, providing researchers with flexible experimental design options.

Societal and environmental impact. OCTAX enables more rigorous evaluation with larger sample sizes, addressing reproducibility concerns that affect institutions with limited computational resources. This implementation can reduce energy consumption compared to resource-intensive benchmarks such as ALE: experiments that once required top-tier clusters can now run efficiently on a single GPU, potentially saving significant compute time and resources.

Limitations. The GPU-based architecture faces performance constraints due to CHIP-8's variable instruction execution complexity. JAX synchronization across parallel environments means each step's execution time depends on the slowest instruction among CHIP-8's 35 operations, typically display rendering or complex ALU operations. The absence of established maximum scores across our game suite prevents the assessment of whether agents approach theoretical performance limits, limiting evaluation of algorithmic performance ceilings.

Future work. OCTAX can expand through community contributions, with hundreds of compatible ROMs available online. The LLM-assisted environment generation pipeline enables curriculum learning and open-ended research through procedurally generated games that provide task diversity. We plan to investigate emulator optimizations including instruction-level parallelization strategies and adaptive batching to address synchronization bottlenecks from variable execution times. We also aim to extend platform support to Super-CHIP8 and XO-CHIP variants: Super-CHIP8 offers higher resolution displays (128×64) and extended instruction sets originally developed for HP48 calculators, while XO-CHIP provides color graphics, improved audio, and expanded memory while maintaining backward compatibility. These extensions would enable OCTAX to support more sophisticated games and visual complexity while preserving the computational efficiency advantages of the JAX-native architecture. Many CHIP-8 games feature multi-agent or multi-player mechanics, which we plan to support in future platform releases. The platform's high-throughput capabilities also position it well for offline RL research, enabling the efficient creation of large-scale datasets and the comprehensive evaluation of offline algorithms across diverse game environments.

REPRODUCIBILITY STATEMENT

We provide complete resources to ensure reproducibility of our results. The OCTAX source code, including all 21 game environment implementations, JAX-based CHIP-8 emulator, and training scripts, is available as supplementary material. Our experimental setup uses standard PPO hyperparameters detailed in Section 4.1, with hardware specifications and performance benchmarking configurations provided in Section 4.2. All training experiments use identical network architectures and hyperparameters across games, enabling direct replication of our learning curves in Figure 3. For the LLM-assisted environment generation pipeline in Section 4.3, we include the prompt templates and generated CHIP-8 assembly code in Appendix E. The modular design of OCTAX allows researchers to extend our game suite using the technical specifications in Section 3. The anonymized repository containing all source code, experiments, and data is available at: https://anonymous.4open.science/r/octax-C8E8/README.md

REFERENCES

Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320, 2021.

- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of artificial intelligence research*, 47: 253–279, 2013.
 - Clément Bonnet, Daniel Luo, Donal Byrne, Shikha Surana, Sasha Abramowitz, Paul Duckworth, Vincent Coyette, Laurence I Midgley, Elshadai Tegegn, Tristan Kalloniatis, et al. Jumanji: a diverse suite of scalable reinforcement learning environments in jax. *arXiv preprint* arXiv:2306.09884, 2023.
 - James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018a. URL http://github.com/jax-ml/jax.
 - James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, et al. Jax: composable transformations of python+ numpy programs. 2018b.
 - Cédric Colas, Olivier Sigaud, and Pierre-Yves Oudeyer. How many random seeds? statistical power analysis in deep reinforcement learning experiments. *arXiv preprint arXiv:1806.08295*, 2018.
 - Steven Dalton et al. Accelerating reinforcement learning through gpu atari emulation. *Advances in Neural Information Processing Systems*, 33:19773–19782, 2020.
 - Sam Earle, Graham Todd, Yuchen Li, Ahmed Khalifa, Muhammad Umair Nasir, Zehua Jiang, Andrzej Banburski-Fahey, and Julian Togelius. Puzzlejax: A benchmark for reasoning and learning. *arXiv* preprint arXiv:2508.16821, 2025.
 - Maxence Faldor, Jenny Zhang, Antoine Cully, and Jeff Clune. Omni-epic: Open-endedness via models of human notions of interestingness with environments programmed in code. *arXiv* preprint *arXiv*:2405.15568, 2024.
 - C Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. Brax–a differentiable physics engine for large scale rigid body simulation. *arXiv preprint arXiv:2106.13281*, 2021.
 - Jason Gardner, Ayan Dutta, Swapnoneel Roy, O Patrick Kreidl, and Ladislau Boloni. Greener deep reinforcement learning: Analysis of energy and carbon efficiency across atari benchmarks. *arXiv* preprint arXiv:2509.05273, 2025.
 - Danijar Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*, 2021.
 - Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
 - Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
 - Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Kinal Mehta, and João G.M. Araújo. Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022. URL http://jmlr.org/papers/v23/21-1342.html.
 - Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv* preprint *arXiv*:1412.6980, 2014.
 - Sotetsu Koyamada, Shinri Okano, Soichiro Nishimori, Yu Murata, Keigo Habara, Haruka Kita, and Shin Ishii. Pgx: Hardware-accelerated parallel game simulators for reinforcement learning. *Advances in Neural Information Processing Systems*, 36:45716–45743, 2023.

- Heinrich Küttler, Nantas Nardelli, Alexander Miller, Roberta Raileanu, Marco Selvatici, Edward
 Grefenstette, and Tim Rocktäschel. The nethack learning environment. Advances in Neural
 Information Processing Systems, 33:7671–7684, 2020.
- Robert Tjarko Lange. gymnax: A JAX-based reinforcement learning environment library, 2022. URL http://github.com/RobertTLange/gymnax.
 - Jarek Liesen, Chris Lu, and Robert Lange. rejax, 2024. URL https://github.com/keraJLi/rejax.
 - Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
 - Timothée Mathieu, Riccardo Della Vecchia, Alena Shilova, Matheus Medeiros Centa, Hector Kohler, Odalric-Ambrym Maillard, and Philippe Preux. Adastop: adaptive statistical testing for sound comparisons of deep rl agents. *arXiv preprint arXiv:2306.10882*, 2023.
 - Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
 - Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, Artem Agarkov, Viacheslav Sinii, and Sergey Kolesnikov. Xland-minigrid: Scalable meta-reinforcement learning environments in jax. *Advances in Neural Information Processing Systems*, 37:43809–43835, 2024.
 - Johan S Obando-Ceron and Pablo Samuel Castro. Revisiting rainbow: Promoting more insightful and inclusive deep reinforcement learning research. *arXiv preprint arXiv:2011.14826*, 2020.
 - Ian Osband, Yotam Doron, Matteo Hessel, John Aslanides, Eren Sezener, Andre Saraiva, Katrina McKinney, Tor Lattimore, Csaba Szepesvari, Satinder Singh, et al. Behaviour suite for reinforcement learning. *arXiv* preprint arXiv:1908.03568, 2019.
 - Eduardo Pignatelli, Jarek Liesen, Robert Tjarko Lange, Chris Lu, Pablo Samuel Castro, and Laura Toni. Navix: Scaling minigrid environments with jax. *arXiv preprint arXiv:2407.19396*, 2024.
 - Nikita Rudin, David Hoeller, Philipp Reist, and Marco Hutter. Learning to walk in minutes using massively parallel deep reinforcement learning. In *Conference on robot learning*, pp. 91–100. PMLR, 2022.
 - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
 - Joseph Suarez. Pufferlib 2.0: Reinforcement learning at 1m steps/s. In *Reinforcement Learning Conference*, 2025.
 - Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 2 edition, 2018.
 - Joseph Weisbecker. An easy programming system. *BYTE*, 3(12):108–122, December 1978. First introduction of CHIP-8 programming language.
 - Jiayi Weng, Min Lin, Shengyi Huang, Bo Liu, Denys Makoviichuk, Viktor Makoviychuk, Zichen Liu, Yufan Song, Ting Luo, Yukun Jiang, et al. Envpool: A highly parallel reinforcement learning environment execution engine. *Advances in Neural Information Processing Systems*, 35:22409–22421, 2022.
 - Kenny Young and Tian Tian. Minatar: An atari-inspired testbed for thorough and reproducible reinforcement learning experiments. *arXiv preprint arXiv:1903.03176*, 2019.

A USE OF LARGE LANGUAGE MODELS

We used Large Language Models in three capacities during this research. First, Claude Opus 4.1 serves as a core research component in Section 4.3, where we demonstrate automated CHIP-8 game generation from high-level descriptions. This represents a novel research contribution, with all generated code validated through compilation and RL experiments. Second, we employed Claude Sonnet 4 for writing assistance, including text refinement, rephrasing technical concepts, and improving academic tone. Third, LLMs generated code documentation, docstrings, and tutorial content. All research ideas, experimental design, and scientific claims originate from the authors. We did not use LLMs for ideation, hypothesis formation, or result interpretation. We manually reviewed and validated all LLM-assisted content for accuracy and take full responsibility for all presented content.

B CHIP-8 TECHNICAL SPECIFICATIONS

B.1 PLATFORM OVERVIEW

CHIP-8 was created by Joseph Weisbecker at RCA in the mid-1970s as a virtual machine for early microcomputers. The platform established one of the first successful portable gaming ecosystems by providing a hardware abstraction layer that enabled games to run across different systems.

B.2 System Architecture

The CHIP-8 architecture consists of:

- Memory: 4KB total, with programs loaded at address 0x200
- **Registers:** 16 8-bit registers (V0-VF), with VF serving as a flag register
- **Display:** 64×32 pixel monochrome screen with XOR-based rendering
- **Input:** 16-key hexadecimal keypad (0-9, A-F)
- Timers: 60Hz delay timer and sound timer
- Audio: Single-tone buzzer

B.3 Instruction Set Highlights

CHIP-8's 35-instruction set includes specialized gaming primitives:

- Sprite Drawing (DXYN): XOR-based rendering enabling collision detection
- Key Input (EX9E, EXA1): Skip instructions based on key state
- BCD Conversion (FX33): Convert register values to decimal display
- Memory Operations: Bulk register loading/storing (FX55, FX65)

The XOR-based sprite system is particularly elegant: drawing the same sprite twice erases it, enabling simple animation and automatic collision detection when pixels turn off.

B.4 FONT SYSTEM

CHIP-8 includes built-in 4×5 pixel font data for hexadecimal digits (0-F), stored at addresses 0x050-0x09F. Games reference these fonts for score and text display by setting the index register to the appropriate font location.

C GAME ENVIRONMENT IMPLEMENTATION DETAILS

C.1 Score Detection Methodology

CHIP-8 games store scoring information in arbitrary memory locations using game-specific formats. Our automated detection operates in two phases:

versus decreasing values (likely lives/health).

energy pills and 2 ghosts

controls and speed fixes

C.2.2 TIMING, PREDICTION & REACTIVE CONTROL

C.2.1 LONG-HORIZON PLANNING & SPATIAL REASONING

Requires strategic thinking, spatial awareness, and multi-step planning

C.2 GAME LIST

dropping

648

649

650

651

652

653 654

655 656

657

658 659

660

661

662

663

664

666

667

669	Requires precise timing, trajectory prediction, and fast reactive responses
670	• brix – Brix by Andreas Gustafsson (1990): Breakout clone with paddle controlling the ball
671	to destroy bricks, 5 lives
672	• pong – Pong: Single player pong game with paddle control
673 674	• squash – Squash by David Winter (1997): Bounce ball around squash court with paddle
675 676	• vertical_brix – Vertical Brix by Paul Robson (1996): Breakout variant with vertical brick layout and paddle movement
677 678	• wipe_off – Wipe Off by Joseph Weisbecker: Move paddle to wipe out spots, 1 point per spot, 20 balls
679 680	• filter – Filter: Catch drops from pipe with paddle
681	C.2.3 RESOURCE MANAGEMENT & TACTICAL DECISIONS
682 683	Requires managing limited resources and making strategic tactical choices
684 685 686	• missile – Missile Command by David Winter (1996): Shoot 8 targets using key 8, earn 5 points per target, 12 missiles total
687 688	 rocket – Rocket by Joseph Weisbecker (1978): Launch rockets to hit moving UFO across top of screen, 9 rockets total
689 690	• submarine – Submarine by Carmelo Cortez (1978): Fire depth charges at submarines, 15 points for small, 5 for large subs
691 692	• tank – Tank Battle: Tank with 25 bombs to hit mobile target, lose 5 bombs if tank hits target
693 694 695	 ufo – UFO by Lutz V (1992): Stationary launcher shoots in 3 directions at flying objects, 15 missiles
696	C.2.4 Exploration & Continuous Navigation
697 698	Requires spatial exploration, obstacle avoidance, and continuous movement control
699 700	• cavern – Cavern by Matthew Mikolay (2014): Navigate cave without hitting walls, modified for leftward exploration
701	• flight_runner – Flight Runner by TodPunk (2014): Simple flight navigation game

Static Analysis: We analyze ROM structure for common programming patterns, particularly binary-

Dynamic Monitoring: During human gameplay sessions, we monitor memory changes to correlate

locations with scoring events. Register trend analysis identifies increasing values (likely scores)

• tetris – Tetris by Fran Dachille (1991): Classic Tetris with piece rotation, movement, and

• blinky – Blinky by Hans Christian Egeberg (1991): Pac-Man clone with 2 lives, maze with

• worm – SuperWorm V4 by RB-Revival Studios (2007): Snake-like game with enhanced

coded decimal (BCD) operations (FX33 instruction) that suggest numeric display routines.

- 702 703
- 704
- 705 706
- 708 709
- 710 711 712
- 713 714 715
- 716 717
- 718 719 720
- 721 722 723
- 724
- 725 726
- 727 728
- 729
- 730 731
- 732 733

736

737 738

739 740 741

742 743 744

745 746

747 748

749 750

751 752 753

754 755

- space_flight Space Flight by Unknown (19xx): Fly through asteroid field using ship navigation controls
- spacejam Spacejam! by William Donnelly (2015): Enhanced ship tunnel navigation game

C.2.5 SIMPLE REACTION & TIMING

Requires basic reaction time and simple decision making

- airplane Airplane: Bombing game where you drop bombs by pressing key 8
- deep Deep8 by John Earnest (2014): Move boat left/right, drop and detonate bombs to destroy incoming squid
- shooting stars Shooting Stars by Philip Baltzer (1978): Classic shooting game

HYPERPARAMETER OPTIMIZATION RESULTS

We conducted a comprehensive grid search on the Pong environment to identify optimal PPO hyperparameters before evaluating across the full game suite. The search explored four key dimensions: number of parallel environments, rollout length, minibatch size, and learning rate. All experiments used 4 epochs per update, GAE lambda of 0.95, and gradient clipping at 0.5.

D.1 SEARCH SPACE

The hyperparameter search explored the following ranges:

- Environments: {128, 256, 512, 1024}
- Rollout steps: {32, 64, 128, 512}
- Minibatches: {4, 8, 16, 32}
- Learning rate: $\{2.5 \times 10^{-4}, 5 \times 10^{-4}, 1 \times 10^{-3}\}$

Each configuration was trained for 1M timesteps with evaluation every 65,536 steps. Final evaluation scores represent the last recorded performance, where less negative values indicate better performance.

D.2 RESULTS SUMMARY

Table 2 presents the key configurations and their final evaluation scores. Higher scores indicate better performance (scores are negative, with values closer to zero being better).

D.3 ANALYSIS AND KEY FINDINGS

Learning rate impact. Higher learning rates significantly improved performance, with 5×10^{-4} and 1×10^{-3} substantially outperforming 2.5×10^{-4} . The top three configurations all used learning rates above the commonly used 2.5×10^{-4} .

Environment scaling. 512 parallel environments provided the optimal balance between computational efficiency and sample diversity. Configurations with 1024 environments showed diminishing returns, possibly due to computational overhead or reduced gradient update frequency.

Rollout length. Shorter rollouts (32 steps) consistently outperformed longer ones, indicating more frequent policy updates may be beneficial for this environment.

Minibatch size. Larger minibatch sizes (16-32) generally improved performance by providing more stable gradient estimates, though the effect was less pronounced than learning rate changes.

Table 2: Hyperparameter search results on Pong environment. Configurations sorted by final evaluation score.

Envs	Steps	Minibatches	LR	Score
512	32	32	0.0005	-2.34
512	32	16	0.001	-2.48
512	32	32	0.001	-2.69
128	128	16	0.00025	-2.95
128	64	8	0.00025	-3.19
512	32	16	0.00025	-3.20
256	64	16	0.00025	-3.38
128	32	4	0.00025	-3.44
128	64	16	0.00025	-3.53
512	32	16	0.0005	-3.73
256	32	4	0.00025	-3.78
128	128	8	0.00025	-3.91
256	128	32	0.00025	-4.03
512	64	16	0.00025	-4.17
1024	32	32	0.00025	-4.34
1024	32	16	0.00025	-4.44
256	128	16	0.00025	-4.66
1024	64	32	0.00025	-4.96

D.4 FINAL CONFIGURATION

Based on these results, we selected the following hyperparameters for all subsequent experiments:

• Parallel environments: 512

Rollout steps: 32Training epochs: 4Minibatches: 32

Learning rate: 5 × 10⁻⁴
GAE lambda: 0.95
Clip epsilon: 0.2

Value function coefficient: 0.5 Entropy coefficient: 0.01

E LLM-Assisted Environment Generation

This appendix details the automated environment generation pipeline using large language models (LLMs) to create novel CHIP-8 games for reinforcement learning research. We demonstrate the complete process from prompt engineering to code generation across three difficulty levels of a Target Shooter game.

E.1 PROMPT ENGINEERING

Our LLM generation pipeline relies on carefully crafted prompts that provide comprehensive CHIP-8 programming context and specific game requirements. The core prompt structure includes CHIP-8 architectural constraints, Octo assembly language syntax, and reinforcement learning compatibility requirements.

Listing 1: LLM prompt template for CHIP-8 game generation

```
You are a **professional CHIP-8 (classic version) game developer**.

Your task is to **design and implement new CHIP-8 games in Octo assembly language**. I will provide you with tutorials and references for Octo assembly. You must be rigorous and ensure that your code is ** syntactically correct, runnable, and follows CHIP-8 conventions**.
```

```
810
811
       <documentation></documentation>
812
813
       <tutorial1></tutorial1>
814
       <tutorial2></tutorial2>
815
816
       <example></example>
817
818
       The goal is to create a **game suitable for reinforcement learning (RL)
           research**, which means:
819
       * The **score** must be stored in a clear and consistent register or
820
           memory location.
821
       * The **termination condition** (game over) must also be easily
822
           extractable (e.g., through a specific flag or register value).
       \star The game should have \star\star \text{deterministic rules} \star\star and be lightweight enough
823
           for training agents.
824
825
       Here is the description of the game you must implement:
826
827
       <description>{{description}}</description>
828
```

The prompt incorporates several key components:

- Role specification: Establishes the LLM as a professional CHIP-8 developer
- Technical constraints: Emphasizes syntactic correctness and CHIP-8 compliance
- RL compatibility: Specifies requirements for score tracking and termination detection
- Reference material: Includes comprehensive CHIP-8 documentation and examples
- Game description: Placeholder for specific game mechanics and objectives

The prompt template includes placeholder tags that are populated with comprehensive CHIP-8 resources: <documentation> contains the official Octo Manual (https://johnearnest.github.io/Octo/docs/Manual.html), <tutorial1> includes the Beginner's Guide (https://johnearnest.github.io/Octo/docs/BeginnersGuide.html), <tutorial2> incorporates the Intermediate Guide (https://johnearnest.github.io/Octo/docs/IntermediateGuide.html), and <example> provides a complete game implementation (https://github.com/JohnEarnest/chip8Archive/blob/master/src/outlaw/outlaw.80) to demonstrate best practices and coding patterns.

E.2 GENERATED TARGET SHOOTER IMPLEMENTATION

Using the prompt template, we generated three progressive difficulty levels of a Target Shooter game. Each level maintains consistent register mappings for score and termination while introducing increasing complexity in target behavior and timing constraints.

E.2.1 LEVEL 1: STATIC TARGETS

The first difficulty level features stationary targets that appear at random locations, focusing on basic aiming and shooting mechanics.

Listing 2: Level 1 Target Shooter - Static targets

```
864
         - E to shoot
865
866
      # Score is stored in register v2 (score_reg)
867
        Game over flag in register v3 (gameover_reg)
868
      # Game ends after hitting 10 targets.
869
870
      871
872
      # Sprite data
      : crosshair
873
             0b10000001
874
              0b01011010
875
              0b00100100
876
              0b01011010
877
              0b01011010
              0b00100100
878
              0b01011010
879
              0b10000001
880
881
      : target
       0b00111100
882
              0b01000010
883
              0h10011001
884
              0b10100101
885
              0b10100101
886
              0b10011001
              0b01000010
887
              0b00111100
888
889
890
      891
      # Register Map - Critical for RL extraction
892
      893
      :alias crosshair_x
                          v0 # Crosshair X position
894
                           v1 # Crosshair Y position
      :alias crosshair_y
                         v2 # SCORE - RL agents read this!
895
      :alias score_reg
896
      :alias gameover_reg v3 # GAME OVER FLAG (0=playing, 1=over)
      :alias target_x v4 # Target X position
:alias target_y v5 # Target Y position
897
898
      :alias target_active v6 # Target active flag
899
      :alias target_active vo # larget active liag
:alias temp1 v7 # Temporary register
:alias temp2 v8 # Temporary register
:alias shot_active v9 # Shot in progress flag
:alias targets_hit va # Count of targets hit (max 10)
:alias key_reg vb # Key input register
900
901
902
903
904
      :const MAX_TARGETS 10
                                # Game ends after 10 targets
905
                                 # Target sprite size
      :const TARGET_SIZE 8
906
                                # Crosshair sprite size
      :const CROSSHAIR_SIZE 8
907
      :const POINTS_PER_HIT 1
                                 # Points awarded per target hit
908
      909
      # Main Game Entry Point
910
      911
912
      : main
        # Initialize game state
913
        score_reg := 0 # Score starts at 0
914
                            # Game is not over
        gameover_reg := 0
915
        targets_hit := 0
                             # No targets hit yet
                            # No target active initially
916
        target_active := 0
        shot_active := 0 # No shot in progress
917
```

```
918
       # Initial crosshair position (center)
919
       crosshair_x := 28
920
       crosshair_y := 12
921
       clear
922
923
       # Draw initial UI
924
       draw-crosshair
925
926
       # Main game loop
927
       loop
         # Check if game should end
928
         if targets_hit == MAX_TARGETS then jump game-over
929
930
         # Spawn new target if none active
         if target_active == 0 then spawn-target
931
932
         # Handle player input
933
         handle-input
934
         # Check for hit if shot was fired
935
         if shot_active == 1 then check-hit
936
937
         # Small delay for playability
938
         temp1 := 1
939
         delay := temp1
940
         wait-delay
941
       again
942
943
      944
      # Game Over Handler
945
      946
      : game-over
947
       gameover_reg := 1
                           # Set game over flag for RL agent
948
949
       # Flash screen to indicate game over
950
       temp1 := 0
       loop
951
         clear
952
         temp2 := 5
953
         delay := temp2
954
         wait-delay
955
         draw-crosshair
956
         if target_active == 1 then draw-target
957
         temp2 := 5
958
         delay := temp2
959
         wait-delay
960
         temp1 += 1
961
         if temp1 != 3 then
962
       again
963
964
       # Infinite loop - game is over
965
         # RL agent should detect gameover_reg == 1
966
       again
967
968
      969
      # Input Handling
970
      971
      : handle-input
```

```
972
        # Save current position
973
        temp1 := crosshair_x
974
        temp2 := crosshair_y
975
        # Movement controls (WASD) - use consistent key codes
976
        key_reg := 7 # A key (left)
977
        if key_reg key then temp1 += -2
978
979
        key_reg := 9 # D key (right)
980
        if key_reg key then temp1 += 2
981
        key\_reg := 5 \# W key (up)
982
        if key_reg key then temp2 += -2
983
984
        key_reg := 8 # S key (down)
        if key_reg key then temp2 += 2
985
986
        # Boundary checking
987
        if temp1 >= 254 then temp1 := 0
                                          # Left boundary (wrapping check)
988
        if temp1 >= 56 then temp1 := 56
                                          # Right boundary
989
        if temp2 >= 254 then temp2 := 0
                                         # Top boundary (wrapping check)
        if temp2 >= 24 then temp2 := 24
990
                                         # Bottom boundary
991
        # Check if position changed
992
        if temp1 != crosshair_x then jump update-crosshair
993
        if temp2 != crosshair_y then jump update-crosshair
994
        # Check for shoot (E key)
995
        key_reg := 6
996
        if key_reg key then shot_active := 1
997
998
        return
999
      : update-crosshair
1000
        # Erase old crosshair
1001
        i := crosshair
1002
        sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1003
1004
        # Update position
        crosshair_x := temp1
1005
        crosshair_y := temp2
1006
1007
        # Draw new crosshair
1008
        sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1009
        # Check shoot after movement
1010
        key_reg := 6
1011
        if key_reg key then shot_active := 1
1012
1013
1014
      # Target Management
1015
      1016
1017
      : spawn-target
1018
        # Generate random position for target
1019
        target_x := random 0x37 # 0-55 range
        target_y := random 0x17 # 0-23 range
1020
1021
        # Ensure minimum distance from edges
1022
        if target_x <= 2 then target_x := 3</pre>
1023
        if target_y <= 2 then target_y := 3</pre>
1024
1025
        target_active := 1
        draw-target
```

```
1026
1027
1028
      : draw-target
1029
        i := target
        sprite target_x target_y TARGET_SIZE
1030
1031
1032
      : draw-crosshair
1033
        i := crosshair
1034
        sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1035
1036
      1037
      # Hit Detection
1038
      1039
      : check-hit
1040
        shot_active := 0 # Reset shot flag
1041
1042
        # Check if target is active
1043
        if target_active == 0 then return
1044
        # Simple hit detection - check if crosshair center is near target
1045
            center
1046
        # Calculate X distance
1047
        temp1 := crosshair_x
1048
        temp1 += 4 # Crosshair center
1049
        temp2 := target_x
        temp2 += 4 # Target center
1050
1051
        # Check X proximity
1052
        if temp1 > temp2 then jump check-x-greater
1053
1054
        \# crosshair is left of \mathbf{or} at target
        temp2 -= temp1
1055
        if temp2 > 6 then return # Too far
1056
        jump check-y-axis
1057
1058
      : check-x-greater
        # crosshair is right of target
1059
        temp1 -= temp2
1060
        if temp1 > 6 then return # Too far
1061
1062
      : check-y-axis
1063
        # Calculate Y distance
        temp1 := crosshair_y
1064
        temp1 += 4 # Crosshair center
1065
        temp2 := target_y
1066
        temp2 += 4 # Target center
1067
1068
        # Check Y proximity
        if temp1 > temp2 then jump check-y-greater
1069
1070
        # crosshair is above or at target
1071
        temp2 -= temp1
1072
        if temp2 > 6 then return # Too far
1073
        jump register-hit
1074
      : check-y-greater
1075
        # crosshair is below target
1076
        temp1 -= temp2
1077
        if temp1 > 6 then return # Too far
1078
1079
      : register-hit
        # Hit confirmed!
```

```
1080
       # Erase target
1081
       draw-target
1082
       target_active := 0
1083
       # Update score (for RL agent)
1084
       score_reg += POINTS_PER_HIT
1085
       targets_hit += 1
1086
1087
       # Sound feedback
1088
       temp1 := 3
       buzzer := temp1
1089
1090
1091
     1092
     # Utility Functions
     1093
1094
     : wait-delay
1095
       loop
1096
        temp1 := delay
1097
        if temp1 != 0 then
1098
       again
1099
```

E.2.2 LEVEL 2: TIME-LIMITED TARGETS

1100 1101

1102 1103

1104

1105

1106

The second level introduces time pressure by making targets disappear after a fixed duration, requiring faster decision-making from RL agents.

Listing 3: Level 2 Target Shooter - Time-limited targets

```
1107
      1108
        Target Shooter - RL Training Game
1109
1110
        A deterministic shooting game designed for
1111
        reinforcement learning research.
1112
1113
        Controls:
        - WASD to move crosshair
1114
        - E to shoot
1115
1116
        Score is stored in register v2 (score_reg)
1117
        Game over flag in register v3 (gameover_reg)
1118
        Game ends after 10 targets (hit or missed).
1119
        Targets disappear after ~3 seconds if not hit.
1120
1121
      1122
1123
      # Sprite data
      : crosshair
1124
             0b10000001
1125
             0b01011010
1126
             0b00100100
1127
             0b01011010
             0b01011010
1128
             0b00100100
1129
             0b01011010
1130
             0b10000001
1131
1132
      : target
             0b00111100
1133
             0b01000010
```

```
1134
              0b10011001
1135
             0b10100101
1136
             0b10100101
1137
             0b10011001
              0b01000010
1138
              0b00111100
1139
1140
1141
      1142
      # Register Map - Critical for RL extraction
      1143
                         v0 # Crosshair X position
      :alias crosshair_x
1145
      :alias crosshair_y
                           v1 # Crosshair Y position
1146
      :alias score_reg
                          v2 # SCORE - RL agents read this!
      :alias gameover_reg v3 # GAME OVER FLAG (0=playing, 1=over)
:alias target_x v4 # Target X position
:alias target_y v5 # Target Y position
1147
1148
1149
      :alias target_active v6 # Target active flag
1150
     :alias temp1
                          v7 # Temporary register
1151
      :alias temp2
                         v8 # Temporary register
      :alias shot_active v9 # Shot in progress flag
1152
      :alias targets_total va # Total targets appeared (max 10)
1153
      :alias key_reg vb # Key input register
:alias target_timer vc # Timer for current target
1154
1155
      :alias missed_targets vd # Count of missed targets
1156
      :const MAX_TARGETS 10
                                # Game ends after 10 targets total
1157
                                # Target sprite size
      :const TARGET_SIZE 8
1158
      :const CROSSHAIR SIZE 8
                              # Crosshair sprite size
1159
      :const POINTS_PER_HIT 1
                               # Points awarded per target hit
1160
      :const TARGET_TIMEOUT 60  # Frames before target disappears (~3 sec at
1161
         20fps)
1162
      1163
      # Main Game Entry Point
1164
      1165
1166
      : main
        # Initialize game state
1167
        score_reg := 0 # Score starts at 0
1168
        gameover_reg := 0 # Game is not over
1169
                            # No targets appeared yet
        targets_total := 0
1170
        missed_targets := 0
                            # No missed targets yet
1171
        target_active := 0
                             # No target active initially
        shot_active
                      := 0
                             # No shot in progress
1172
        target_timer := 0
                            # Timer at 0
1173
1174
        # Initial crosshair position (center)
1175
        crosshair_x := 28
1176
        crosshair_y := 12
1177
1178
1179
        # Draw initial UI
1180
        draw-crosshair
1181
        # Main game loop
1182
        loop
1183
          # Check if game should end (10 total targets)
1184
          if targets_total == MAX_TARGETS then jump game-over
1185
1186
          # Spawn new target if none active
          if target_active == 0 then spawn-target
1187
```

```
1188
         # Check target timeout
1189
         if target_active == 1 then check-target-timeout
1190
1191
         # Handle player input
         handle-input
1192
1193
         # Check for hit if shot was fired
1194
         if shot_active == 1 then check-hit
1195
1196
         # Small delay for playability
         temp1 := 1
1197
         delay := temp1
1198
         wait-delay
1199
1200
       again
1201
      1202
      # Target Timeout Check
1203
      1204
1205
      : check-target-timeout
       # Decrement timer
1206
       target\_timer += -1
1207
1208
       \# Check if timer expired
1209
       if target_timer != 0 then return
1210
       # Target timed out - count as miss
1211
       draw-target # Erase target
1212
       target_active := 0
1213
       missed_targets += 1
1214
1215
       # Brief sound to indicate miss
1216
       temp1 := 1
       buzzer := temp1
1217
1218
1219
      1220
      # Game Over Handler
      1221
1222
      : game-over
1223
                           # Set game over flag for RL agent
       gameover_reg := 1
1224
1225
       # Flash screen to indicate game over
       temp1 := 0
1226
       loop
1227
         clear
1228
         temp2 := 5
1229
         delay := temp2
1230
         wait-delay
1231
         draw-crosshair
1232
         if target_active == 1 then draw-target
1233
         temp2 := 5
1234
         delay := temp2
1235
         wait-delay
1236
         temp1 += 1
1237
         if temp1 != 3 then
1238
       again
1239
1240
       # Infinite loop - game is over
1241
       loop
         # RL agent should detect gameover_reg == 1
```

```
1242
        again
1243
1244
      1245
      # Input Handling
      1246
1247
      : handle-input
1248
        # Save current position
1249
       temp1 := crosshair_x
1250
        temp2 := crosshair_y
1251
        # Movement controls (WASD) - use consistent key codes
1252
       key_req := 7 # A key (left)
1253
       if key_reg key then temp1 += -2
1254
       key_reg := 9 # D key (right)
1255
       if key_reg key then temp1 += 2
1256
1257
        key\_reg := 5 \# W key (up)
1258
       if key_reg key then temp2 += -2
1259
       key_reg := 8 # S key (down)
1260
        if key_reg key then temp2 += 2
1261
1262
        # Boundary checking
1263
        if temp1 >= 254 then temp1 := 0
                                       # Left boundary (wrapping check)
1264
        if temp1 >= 56 then temp1 := 56
                                       # Right boundary
       if temp2 >= 254 then temp2 := 0
1265
                                       # Top boundary (wrapping check)
                                      # Bottom boundary
       if temp2 >= 24 then temp2 := 24
1266
1267
        # Check if position changed
1268
        if temp1 != crosshair_x then jump update-crosshair
1269
       if temp2 != crosshair_y then jump update-crosshair
1270
        # Check for shoot (E key)
1271
        key_reg := 6
1272
       if key_reg key then shot_active := 1
1273
1274
       return
1275
      : update-crosshair
1276
        # Erase old crosshair
1277
       i := crosshair
1278
       sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1279
        # Update position
1280
       crosshair_x := temp1
1281
       crosshair_y := temp2
1282
1283
        # Draw new crosshair
1284
       sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1285
        # Check shoot after movement
1286
       key_reg := 6
1287
       if key_reg key then shot_active := 1
1288
1289
      1290
      # Target Management
1291
      1292
1293
      : spawn-target
1294
        # Generate random position for target
       target_x := random 0x37 # 0-55 range
1295
       target_y := random 0x17 # 0-23 range
```

```
1296
1297
        # Ensure minimum distance from edges
1298
        if target_x <= 2 then target_x := 3</pre>
        if target_y <= 2 then target_y := 3</pre>
1299
1300
        target_active := 1
1301
        target_timer := TARGET_TIMEOUT # Set timeout timer
1302
        targets_total += 1
                                        # Increment total targets count
1303
        draw-target
1304
1305
      : draw-target
1306
        i := target
1307
        sprite target_x target_y TARGET_SIZE
1308
1309
      : draw-crosshair
1310
        i := crosshair
1311
        sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1312
1313
      1314
      # Hit Detection
1315
      1316
1317
      : check-hit
1318
        shot_active := 0 # Reset shot flag
1319
        # Check if target is active
1320
        if target_active == 0 then return
1321
1322
        # Simple hit detection - check if crosshair center is near target
1323
            center
        # Calculate X distance
1324
        temp1 := crosshair_x
1325
        temp1 += 4 # Crosshair center
1326
        temp2 := target_x
1327
        temp2 += 4 # Target center
1328
        # Check X proximity
1329
        if temp1 > temp2 then jump check-x-greater
1330
1331
        # crosshair is left of or at target
1332
        temp2 -= temp1
1333
        if temp2 > 6 then return # Too far
        jump check-y-axis
1334
1335
      : check-x-greater
1336
        # crosshair is right of target
1337
        temp1 -= temp2
1338
        if temp1 > 6 then return # Too far
1339
      : check-y-axis
1340
        # Calculate Y distance
1341
        temp1 := crosshair_y
1342
        temp1 += 4 # Crosshair center
1343
        temp2 := target_y
        temp2 += 4 # Target center
1344
1345
        # Check Y proximity
1346
        {\tt if} temp1 > temp2 then jump check-y-greater
1347
1348
        \# crosshair is above \mathbf{or} at target
        temp2 -= temp1
1349
        if temp2 > 6 then return # Too far
```

```
1350
       jump register-hit
1351
1352
      : check-y-greater
1353
       # crosshair is below target
       temp1 -= temp2
1354
       if temp1 > 6 then return # Too far
1355
1356
      : register-hit
1357
       # Hit confirmed!
1358
       # Erase target
       draw-target
1359
       target_active := 0
1360
       target_timer := 0 # Clear timer
1361
1362
       # Update score (for RL agent)
       score_reg += POINTS_PER_HIT
1363
1364
       # Sound feedback
1365
       temp1 := 3
1366
       buzzer := temp1
1367
1368
      1369
      # Utility Functions
1370
      1371
1372
      : wait-delay
1373
       loop
         temp1 := delay
1374
         if temp1 != 0 then
1375
       again
1376
1377
```

E.2.3 Level 3: Moving Targets with Time Constraints

13781379

1380

1381

1382 1383 The most challenging level combines target movement with time limits, requiring predictive aiming and rapid response times.

Listing 4: Level 3 Target Shooter - Moving targets with time constraints

```
1384
     1385
1386
        Target Shooter Level 3 - RL Training Game
1387
1388
       A deterministic shooting game designed for
       reinforcement learning research.
1389
1390
       Controls:
1391
       - WASD to move crosshair
1392
        - E to shoot
1393
       Score is stored in register v2 (score reg)
1394
       Game over flag in register v3 (gameover_reg)
1395
1396
       Features:
1397
       - Moving targets that bounce off walls
       - Targets disappear after ~3 seconds if not hit
1398
       - Game ends after 10 targets (hit or missed)
1399
1400
     1401
1402
     # Sprite data
1403
     : crosshair
            0b10000001
```

```
1404
              0b01011010
1405
              0b00100100
1406
              0b01011010
1407
              0b01011010
              0b00100100
1408
              0b01011010
1409
              0b10000001
1410
1411
      : target
1412
              0b00111100
              0b01000010
1413
              0b10011001
1414
              0b10100101
1415
              0b10100101
1416
              0b10011001
              0b01000010
1417
              0b00111100
1418
1419
1420
      1421
      # Register Map - Critical for RL extraction
      1422
1423
      :alias crosshair_x v0 # Crosshair X position
1424
      :alias crosshair_y v1 # Crosshair Y position
:alias score_reg v2 # SCORE - RL agents read this!
:alias gameover_reg v3 # GAME OVER FLAG (0=playing, 1=over)
1425
1426
      v4 # Target X position
1427
1428
      :alias target_active v6 # Target active flag
1429
      :alias temp1 v7 # Temporary register
:alias temp2 v8 # Temporary register
1430
                          v8 # Temporary register
1431
      :alias shot_active v9 # Shot in progress flag
      :alias targets_total va # Total targets appeared (max 10)
:alias key_reg vb # Key input register
:alias target_timer vc # Timer for current target
1432
1433
1434
                           vd # Target X velocity
      :alias target_vx vd # Target X velocity :alias target_vy ve # Target Y velocity
1435
1436
      :const MAX_TARGETS 10
                                # Game ends after 10 targets total
1437
      :const TARGET_SIZE 8
                                # Target sprite size
1438
      :const CROSSHAIR_SIZE 8  # Crosshair sprite size
1439
      :const POINTS_PER_HIT 1
                                 # Points awarded per target hit
1440
                                # Frames before target disappears (~4 sec with
      :const TARGET_TIMEOUT 80
1441
           movement)
1442
      1443
      # Main Game Entry Point
1444
      1445
1446
      : main
       # Initialize game state
1447
                   := 0  # Score starts at 0
        score_req
1448
        gameover_reg
                       := 0
                              # Game is not over
1449
        targets_total := 0 # No targets appeared yet
1450
       target_active := 0 # No target active initially
1451
       shot_active := 0 # No shot in progress
       target_timer := 0 # Timer at 0
1452
       target_vx := 0 # No initial velocity
1453
                      := 0
                             # No initial velocity
        target_vy
1454
1455
        # Initial crosshair position (center)
1456
        crosshair_x := 28
1457
        crosshair_y := 12
```

```
1458
        clear
1459
1460
        # Draw initial UI
1461
        draw-crosshair
1462
        # Main game loop
1463
        loop
1464
          # Check if game should end (10 total targets)
1465
          if targets_total == MAX_TARGETS then jump game-over
1466
          # Spawn new target if none active
1467
          if target_active == 0 then spawn-target
1468
1469
          # Update target position if active
1470
          if target_active == 1 then move-target
1471
          # Check target timeout
1472
          if target_active == 1 then check-target-timeout
1473
1474
          # Handle player input
1475
          handle-input
1476
          # Check for hit if shot was fired
1477
          if shot_active == 1 then check-hit
1478
1479
          # Small delay for playability
1480
          temp1 := 1
          delay := temp1
1481
          wait-delay
1482
1483
        again
1484
1485
      1486
      # Target Movement
      1487
1488
      : move-target
1489
        # Erase target at current position
1490
        draw-target
1491
        # Update X position
1492
        target_x += target_vx
1493
1494
        # Check X boundaries and bounce
1495
        if target_x >= 250 then jump bounce-left
                                                    # Hit left edge
        if target_x >= 56 then jump bounce-right
                                                   # Hit right edge
1496
1497
      : check-y-movement
1498
        # Update Y position
1499
        target_y += target_vy
1500
        # Check Y boundaries and bounce
1501
        if target_y >= 250 then jump bounce-top
                                                     # Hit top edge
1502
        if target_y >= 24 then jump bounce-bottom
                                                     # Hit bottom edge
1503
1504
      : finish-move
1505
        # Draw target at new position
        draw-target
1506
        return
1507
1508
      : bounce-left
1509
        target_x := 1
1510
        target_vx := 1 # Reverse to move right
        jump check-y-movement
1511
```

```
1512
      : bounce-right
1513
       target_x := 55
1514
       target_vx := 255 # -1 to move left
1515
       jump check-y-movement
1516
      : bounce-top
1517
       target_y := 1
1518
       target_vy := 1 # Reverse to move down
1519
       jump finish-move
1520
      : bounce-bottom
1521
       target_y := 23
1522
       target_vy := 255
                       \# -1 to move up
1523
       jump finish-move
1524
      1525
      # Target Timeout Check
1526
      1527
1528
      : check-target-timeout
1529
       # Decrement timer
       target\_timer += -1
1530
1531
       # Check if timer expired
1532
       if target_timer != 0 then return
1533
1534
       # Target timed out - count as miss
       draw-target # Erase target
1535
       target_active := 0
1536
1537
       # Brief sound to indicate miss
1538
       temp1 := 1
1539
       buzzer := temp1
1540
1541
      1542
      # Game Over Handler
1543
      1544
1545
      : game-over
       gameover_reg := 1
                           # Set game over flag for RL agent
1546
1547
       # Flash screen to indicate game over
1548
       temp1 := 0
1549
       loop
1550
         clear
         temp2 := 5
1551
         delay := temp2
1552
         wait-delay
1553
1554
         draw-crosshair
         if target_active == 1 then draw-target
1555
         temp2 := 5
1556
         delay := temp2
1557
         wait-delay
1558
1559
         temp1 += 1
         if temp1 != 3 then
1560
       again
1561
1562
       # Infinite loop - game is over
1563
       loop
1564
         # RL agent should detect gameover_reg == 1
1565
       again
```

```
1566
      1567
      # Input Handling
1568
      1569
      : handle-input
1570
        # Save current position
1571
        temp1 := crosshair_x
1572
        temp2 := crosshair_y
1573
1574
        # Movement controls (WASD) - use consistent key codes
        key_reg := 7 # A key (left)
1575
        if key_reg key then temp1 += -2
1576
1577
        key_reg := 9 # D key (right)
1578
        if key_reg key then temp1 += 2
1579
        key_reg := 5  # W key (up)
1580
        if key_reg key then temp2 += -2
1581
1582
        key_reg := 8 # S key (down)
1583
        if key_reg key then temp2 += 2
1584
        # Boundary checking
1585
        if temp1 >= 254 then temp1 := 0
                                       # Left boundary (wrapping check)
1586
                                       # Right boundary
        if temp1 >= 56 then temp1 := 56
1587
                                       # Top boundary (wrapping check)
        if temp2 \Rightarrow 254 then temp2 \Rightarrow 0
1588
        if temp2 >= 24 then temp2 := 24
                                       # Bottom boundary
1589
        # Check if position changed
1590
        if temp1 != crosshair_x then jump update-crosshair
        if temp2 != crosshair_y then jump update-crosshair
1592
1593
        # Check for shoot (E key)
        key_reg := 6
1594
        if key_reg key then shot_active := 1
1595
1596
        return
1597
1598
      : update-crosshair
        # Erase old crosshair
1599
        i := crosshair
1600
        sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1601
1602
        # Update position
1603
        crosshair_x := temp1
        crosshair_y := temp2
1604
1605
        # Draw new crosshair
1606
        sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1607
1608
        # Check shoot after movement
        key_reg := 6
1609
        if key_reg key then shot_active := 1
1610
1611
1612
      1613
      # Target Management
      1614
1615
      : spawn-target
1616
        # Generate random position for target
1617
       target_x := random 0x37 # 0-55 range
        target_y := random 0x17 # 0-23 range
1618
1619
        # Ensure minimum distance from edges
```

```
1620
        if target_x <= 2 then target_x := 3</pre>
1621
        if target_y <= 2 then target_y := 3</pre>
1622
1623
        # Generate random velocity (-1, 0, or 1 for each axis)
        target_vx := random 0x03
1624
        if target_vx == 2 then target_vx := 255 \# Convert 2 to -1
1625
1626
        target_vy := random 0x03
1627
        if target_vy == 2 then target_vy := 255 \# Convert 2 to -1
1628
        # Ensure target is moving (not both velocities zero)
1629
        {\tt if} target_vx == 0 then jump ensure-movement
1630
        jump finish-spawn
1631
1632
      : ensure-movement
        if target_vy == 0 then target_vy := 1
1633
1634
      : finish-spawn
1635
        target_active := 1
1636
        target_timer := TARGET_TIMEOUT # Set timeout timer
1637
        targets_total += 1
                                         # Increment total targets count
1638
        draw-target
1639
1640
      : draw-target
1641
        i := target
1642
        sprite target_x target_y TARGET_SIZE
1643
1644
      : draw-crosshair
1645
        i := crosshair
1646
        sprite crosshair_x crosshair_y CROSSHAIR_SIZE
1647
1648
      1649
      # Hit Detection
1650
      1651
1652
      : check-hit
        shot_active := 0 # Reset shot flag
1653
1654
        # Check if target is active
1655
        if target_active == 0 then return
1656
1657
        # Simple hit detection - check if crosshair center is near target
1658
            center
        # Calculate X distance
1659
        temp1 := crosshair_x
1660
        temp1 += 4 # Crosshair center
1661
        temp2 := target_x
1662
        temp2 += 4 # Target center
1663
        # Check X proximity
1664
        if temp1 > temp2 then jump check-x-greater
1665
1666
        # crosshair is left of or at target
1667
        temp2 -= temp1
        if temp2 > 6 then return # Too far
1668
        jump check-y-axis
1669
1670
      : check-x-greater
1671
        # crosshair is right of target
1672
        temp1 -= temp2
        {\tt if} temp1 > 6 then return # Too {\tt far}
1673
```

```
1674
      : check-y-axis
1675
        # Calculate Y distance
1676
        temp1 := crosshair_y
        temp1 += 4 # Crosshair center
1677
        temp2 := target_y
1678
        temp2 += 4 # Target center
1679
1680
        # Check Y proximity
1681
        if temp1 > temp2 then jump check-y-greater
1682
        # crosshair is above or at target
1683
        temp2 -= temp1
1684
        if temp2 > 6 then return # Too far
1685
        jump register-hit
1686
      : check-y-greater
1687
        # crosshair is below target
1688
        temp1 -= temp2
1689
        if temp1 > 6 then return # Too far
1690
1691
      : register-hit
        # Hit confirmed!
1692
        # Erase target
1693
        draw-target
1694
        target_active := 0
1695
        target_timer := 0 # Clear timer
1696
        # Update score (for RL agent)
        score_reg += POINTS_PER_HIT
1698
1699
        # Sound feedback
1700
        temp1 := 3
1701
        buzzer := temp1
1702
1703
      1704
      # Utility Functions
1705
      1706
1707
      : wait-delay
        loop
1708
          temp1 := delay
1709
          if temp1 != 0 then
1710
        again
1711
1712
```

E.2.4 Environment Integration and Wrapper Implementation

1713

1714 1715

1716

1717

1718

1719

1720

1721

1722

1723 1724

1725

1726 1727 Once the LLM generates the CHIP-8 assembly code for each difficulty level, the games require integration with OCTAX's reinforcement learning interface. The environment wrapper extracts reward signals and termination conditions from the consistent register mapping established during code generation.

The Target Shooter implementation demonstrates the integration between LLM-generated content and the OCTAX framework. Each level maintains identical register assignments to ensure compatibility across the difficulty progression, enabling curriculum learning experiments without code modifications.

Listing 5: Target Shooter environment wrapper implementation

```
from octax import EmulatorState

def score_fn(state: EmulatorState) -> float:
    """
```

```
1728
           Extract score from register V[2]
1729
           Score increments by 1 for each successful hit
1730
           Range: 0-10 points
1731
           return state.V[2]
1732
1733
      def terminated_fn(state: EmulatorState) -> bool:
1734
1735
           Check game termination flag in register V[3]
           Game ends after 10 total targets (hit or missed in levels 2-3)
1736
1737
           return state.V[3] == 1
1738
1739
       # CHIP-8 key mapping for controls
1740
       \# W=5 (up), A=7 (left), S=8 (down), D=9 (right), E=6 (shoot)
      action\_set = [5, 7, 8, 9, 6]
1741
1742
      metadata = {
1743
           "title": "Target_Shooter_-_LLM-Generated_RL_Environment",
1744
           "authors": ["Fully_LLM-Generated_Environment"],
1745
           "description": "AI-generated_progressive_difficulty_environment",  
           "roms": {
1746
               "target_shooter_level1": {
1747
                   "file": "target_shooter_level1.ch8",
1748
                   "description": "Static_targets_-_Basic_aiming_skills"
1749
1750
               "target_shooter_level2": {
                    "file": "target_shooter_level2.ch8",
1751
                   "description": "Time-limited_static_targets"
1752
1753
               "target_shooter_level3": {
1754
                   "file": "target_shooter_level3.ch8",
1755
                   "description": "Moving_time-limited_targets"
1756
               }
           }
1757
1758
```

The consistent register mapping across all three levels enables direct comparison of agent performance and facilitates automated curriculum progression. Register V[2] consistently stores the score for reward calculation, while V[3] serves as the binary termination flag. The five-action control scheme (WASD movement plus shoot) provides sufficient complexity for interesting policies while remaining tractable for systematic analysis.