

Adaptive Multiple AI by Leveraging Reinforcement Learning

Ravisankar Popuri, Gunamani Jena, Shubhashish Jena, Chandra Mouli VSA, P Devabalan

Research scholar, CMJ University, Jorabat Meghalaya, Roland Institute of Technology, BPUT

ravi.shinsou@gmail.com, drgjena@gmail.com

Article Received: 10 Dec 2024,

Revised: 16 Jan 2025,

Accepted: 26 Jan 2025

Abstract: This paper is about utilizing reinforcement learning (RL) for the development of adaptive AI systems utilizing Unreal Engine 5. This paper focuses on how reinforcement learning can be used to create more efficient artificial intelligence or Non- Playable Characters (NPCs) that dynamically adjust their behavior based on any given in-game context. This provides not only more immersive and challenging gameplay, but also an easy way for AI to mimic human behavior. The research has a detailed analysis of integrating RL algorithm with unreal engine's blueprint and C++ systems. This provides a framework for developers to implement adaptive AI that can evolve over time. The overall performance of this AI system is evaluated through various gameplay scenarios, that demonstrate significant improvements in player engagement. In games now a days, multiplayer experiences continue to grow in complexity and popularity. Game developers are seeking more adaptive and strategically efficient Non-Playable Characters (NPCs) to mimic human-like behavior. This paper showcases an efficient AI system where two opposing teams of agents which are trained using Proximal Policy Optimization (PPO) engage in open-field combat within a game level under Unreal Engine. Each agent aims to maximize individual and collective team rewards by eliminating opponent team members. Through iterative training and carefully designed reward structures, this AI system learns robust combat strategies which leads to an efficient combat AI that exhibits sophisticated coordination behaviors. Results indicate that PPO-driven agents can adapt to various geometrical game levels more effectively compared to conventional scripted NPCs or enemy AI. This paper demonstrates the potential of deep reinforcement learning for complex multi-agent in diverse game environments.

Keywords: create Non- AI Behavior, Playable Characters (NPCs), Proximal Policy Optimization (PPO), RNN-based PPO, Simultaneous multi-RL agents.

I. INTRODUCTION

Game AI has traditionally relied on hard coded decision trees or finite-state machines to drive / control Non-Playable Characters (NPCs). These methods can produce engaging experiences, but they often fail to adapt efficiently to diverse game environments. Reinforcement Learning (RL) presents a proficient alternative by combining deep neural networks. AI agents are trained through repeated interactions in game level environment, RL has a huge potential to implement AI behavior strategies that outperform scripted methods in order to dynamically adapt to changing conditions. This study focuses on designing a multi-agent AI system with two teams, comprised of autonomous agents. These agents are trained using a specific RL algorithm called Proximal Policy Optimization (PPO). All AI agents will face off in a simple open world map created in Unreal Engine 5. Each agent's aim is to shoot its opponent team's agents to contribute to the team's overall success. PPO is mainly selected for this purpose due to its balance between sample efficiency and stable policy updates, attributes. These factors make PPO a strong candidate for multi-agent tasks with continuous or semi-continuous action spaces.

I.I. Contribution

Multi-agent PPO powered AI agents operate as two opposing teams. An in-depth and hands-on implementation of PPO algorithm that accommodates advanced dynamics and complex 3D navigation. Quantitative and qualitative evaluations highlight how RL-driven agents adapt to evolving circumstances. RL agents demonstrate an efficient coordinated team strategy which mimics multiplayer behavior although being single player. Optimizable and extendable framework that can evolve for future needs.

II. LITERATURE REVIEW

PPO was introduced by Schulman et al. (2017), is a policy gradient method designed to stabilize the training of neural-network-based RL agents. PPO modifies the objective function in such a way that large, destabilizing policy updates are penalized further ensuring smoother convergence. In single-agent settings, PPO has been applied successfully to continuous control tasks, arcade-style video games, and simulated robotics. In multi-agent

scenarios, it retains its stability advantages but requires careful design of shared reward structures and training loops.

Traditional RL algorithms like Q-learning approximate the value of actions (Q-values) and select actions that maximize expected returns. In contrast, policy gradient methods directly learn the probability distribution of actions given states, typically represented by a parameterized function $\pi_{\theta}(a | s)$, where θ are the parameters (e.g., weights in a neural network). The agent's objective is to maximize the expected return, often formalized as:

$$J(\theta) = E_{\tau \sim \pi_{\theta}} [r = \sum_{t=0}^T \gamma^t R_t],$$

where τ is a trajectory (sequence of states, actions, and rewards), T is the episode length, γ is the discount factor, and R_t is the reward at time t .

One challenge with policy gradient methods is the risk of destructive large updates, which can push the policy far from previously successful behaviors and destabilize training. PPO addresses this by introducing a clipped objective that penalizes deviations from the old policy beyond a certain threshold. This ensures iterative, proximal updates, striking a balance between exploring new strategies and consolidating previous gains.

Steps involved in PPO algorithm:

- *Data Collection:* The current policy π_{θ} is used to run multiple episodes or mini episodes, generating a batch of trajectories.
- *Advantage Estimation:* An advantage function A_t is calculated for each timestep t , indicating how much better taking action a_t was compared to the agent's baseline expectation.
- *Policy Update:* The ratio $r_t(\theta)$ between the new policy and old policy probabilities for each action is calculated: $r_t(\theta) = \pi_{\theta_{\text{old}}}(a_t | s_t) / \pi_{\theta}(a_t | s_t)$.
- Then, the clipped objective: $L_{\text{CLIP}}(\theta) = E_t[\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)]$, is optimized, where ϵ (commonly 0.1 or 0.2) controls the clipping range.
- *Value Function Update:* PPO often also learns a value function $V_{\phi}(s_t)$ to aid in advantage estimation, updating ϕ via a mean-squared error (MSE) loss.
- *Repeat:* New experiences are collected with the updated policy, and the steps repeat until convergence or for a set number of iterations.

By constraining each policy update, PPO yields a learning process characterized by greater stability and fewer catastrophic policy shifts. This characteristic is especially advantageous in multi-agent scenarios where each agent's policy update can affect the entire environment dynamics.

[Mastering the game of Go with deep neural networks and tree search] demonstrates a unique approach to playing board game Go at very high superhuman level by combining deep neural networks with a Monte Carlo Tree Search (MCTS) framework. A policy network is trained to predict human strategies, while a value network evaluates board positions. The system adapts its strategy through self-play. This sets a new precedent for the application of reinforcement learning in game environments.

[Grandmaster level in StarCraft II using multi-agent reinforcement learning] explores a multi-agent reinforcement learning approach for the real-time strategy game StarCraft II. This utilizes a combination of deep learning architectures, distributed training and centralized value functions to coordinate multiple in-game units. The findings illustrate the scalability of RL methods and their ability to handle the complexity of strategy games.

III. PROPOSED APPROACH

III.I. Game AI Using Unreal Engine

Unreal Engine is a very popular game engine that is widely adopted in both academic research and commercial game development. It has high-fidelity graphics, physics simulations, and flexible scripting environments. Game developers and researchers are using it for studying emergent behaviors in AI-driven agents, pathfinding in

complex terrains, and player adaptation. Unreal engine can be regarded as an IDE where the entire system can be implemented without depending on any other external tools.

In the proposed system two teams governed by reinforced learning agents will be trained to shoot multiple targets by utilizing RL agent interactor, RL agent Trainer and RL agent Manager with a 3D game character. Interactor serves as the hands and eyes of the RL agents. The trainer handles the core learning algorithm, storing data, updating models, and providing the policy decisions. Finally, the Manager orchestrates the entire simulation by instantiating agents.

III.II. Architecture of PPO

The PPO Architecture diagram depicted in Fig.1 showcases the major components of the multi-agent PPO system in Unreal Engine. The diagram illustrates the data flow between the modules involved along with their interaction during training and inference.

Component Explanation

- i. **Parameter Server:** Parameter Server is a central coordinator which maintains a centralized repository of the shared policy ($\pi\theta$) and value function networks which serves as the “authority” on which weights each agent should use at any given point of time.

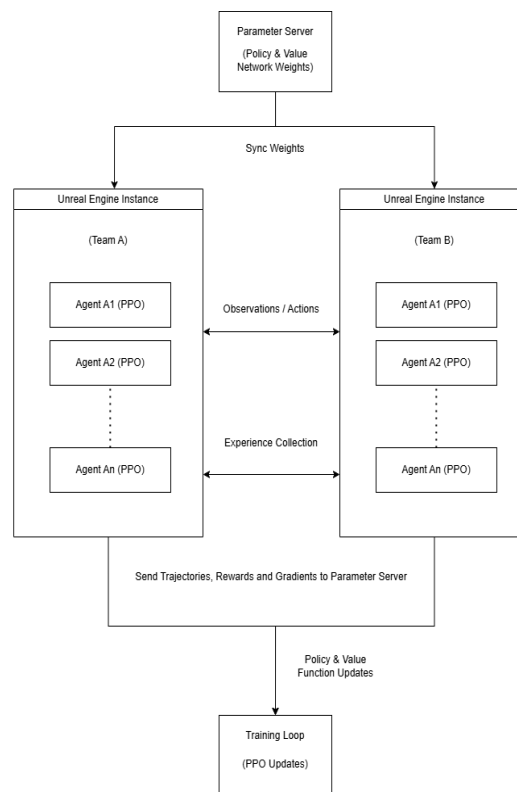
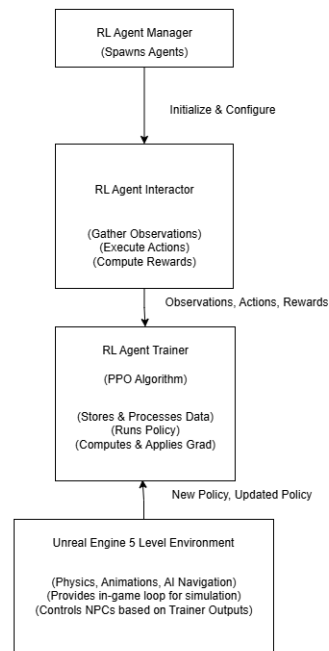


Fig. 1(a) PPO Architecture

Fig. 1(b) Unreal Engine RL Architecture

- ii. **Data Flow:** Sync Weights periodically distributes the latest network parameters (θ for the policy network and ϕ for the value network) to each Team instance.
- iii. **Policy & Value Function Updates:** After receiving experience data from the agents like states, actions, rewards and next states, it performs PPO training steps like policy gradient updates, clipping, advantage calculation. Then the updates weights are synchronized back to the agents.



PPO Agents

Local Decision Making: Every agent $A_1, B_1 \dots A_n, B_n$ observes the game states like health, positions, line of sight to decide on actions like movement, firing and reposition in accordance with the opponents on the field using the local copy of the policy network.

- i. **Experience Collection:** At each timestep the agent logs a tuple $(st, at, rt, st + 1)$. These are assembled into trajectories and then sent back to the Parameter Server for training.
- ii. **Actions & Observations:** Agents receive data like positions of enemies, team status from the Unreal Engine environment and output discrete or continuous actions like move forward, fire weapon.
- iii. **Experience Collection:** Agents gather rewards like kill bonuses like added score, death penalty and store them in memory until they reach a designated batch size of the end of an episode.
- iv. **Training Loop:** When an episode or mini- batch is completed (e.g. iterations / steps), the local agents send their collected trajectories to the parameter server.
- v. **Policy & Value Network Optimization:** The server applies PPO updates to compute the advantage function using the value network, calculate the clipped objective to ensure stable updates and finally adjust network parameters θ (policy) and ϕ (value function) via gradient descent.
- vi. **Updated Policy Broadcast:** After each update, new network weights are broadcast to the agents so that they continue acting with the latest policy.

Communication Protocol

Will be realized with a client-server model over sockets or Remote Procedural Call (RPC), ensuring data from multiple Unreal Engine instances flows to a single Parameter Server.

Synchronization Frequency is determined by the training protocol where agents might sync after every episode, every N steps or after a fixed time interval.

Implementation in Unreal Engine

Environment Design

An open-field map is designed in Unreal Engine, populated by two teams. Each team is composed of multiple RL agents. Key considerations include:

- **Terrain Layout:** Mostly flat with some scattered obstacles or cover points, ensuring line-of-sight constraints and possibilities for strategic maneuvering.
- **Spawn Points:** Both teams spawn at opposite ends, forcing them to traverse the field to encounter enemies.
- **Collision and Damage Systems:** Leveraging Unreal Engine’s physics and collision detection for bullet projectiles, hit registration, and health calculations.

State, Action, and Reward Spaces

State Representation

Each agent’s observation is encoded in vectors that capture:

- **Self-Information:** Current health, ammunition count, cooldown timers, movement speed.
- **Teammate Data:** Relative positions, health status of nearby teammates.
- **Enemy Data:** Approximate direction, distance, and known positions of visible enemies.
- **Environment Context:** Spatial position, presence of cover, and navigability of terrain. This data is periodically sampled and fed into a neural network representing $\pi\theta$.

Action Space

Actions are discretized into commands such as:

- **Movement:** Move forward, backward, strafe left/right, sprint. Rotation: Look left/right, adjust pitch for aiming.
- **Combat:** Fire weapon, reload, relocate.
- In some variants, these actions can be combined or subdivided for finer control, potentially moving toward a continuous action space if required.

Reward Structure

We designed a reward function to promote strategic team combat:

- **Kill Reward (+r1):** Positive reward for successfully eliminating an opponent. Team Support Bonus (+r2): Additional reward for assisting a teammate in a kill (e.g., damage contribution). Death Penalty (-r3): Negative reward if the agent is killed. Friendly Fire Penalty (-r4): Discourages shooting teammates.
- **Survival Incentive (+r5):** Minor positive reward for staying alive each timestep to emphasize survival as a contributing factor to victory. The magnitudes of these rewards (r1, r2, r3, r4, r5) are crucial to shaping the emergent policy and are subject to iterative tuning.

IV. RESULTS

Quantitative Analysis

Team Victory Rate dictates the percentage of matches won by each team over 100 evaluation games. Over training time, both teams improved, but the team with better hyperparameter tuning or more stable training outperformed the other.

Mean Episode Return dictates the average cumulative reward per agent per episode. This metric helps gauge individual agent effectiveness.

Kill/Death Ratio demonstrates the offensive prowess and survivability of the agents.

Coordination Score is a custom metric based on proximity and timing of attacks, measuring how often agents worked in unison. Over approximately 3000 training iterations, the PPO agents showed a substantial increase in coordinated behaviors. Many matches culminated in agents successfully flanking enemy positions, suggesting a learned tactical approach rather than random engagements.

Qualitative Observations

Adaptive Positioning where agents began using the terrain strategically, seeking cover in small dips in the ground or behind scattered props. Over time, the friendly-fire penalty effectively diminished instances of teammates firing in each other’s line-of-sight.

In post-evaluation tests, the PPO agents demonstrated a challenging skill ceiling, sometimes catching opponents off-guard with concentrated fire.

Result and Discussion

Fig.2 showcases all agents for both teams getting instantiated / spawned by the RL Agent Manager.

```
LogOnline: OSS: Created online subsystem instance for: :Context_1
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1767430095 with id 0.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1760022275 with id 1.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1793322284 with id 0.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1787741281 with id 1.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1764061277 with id 2.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1792659283 with id 2.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1767242278 with id 3.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1752002279 with id 3.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1752662272 with id 4.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1806128286 with id 4.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1761668276 with id 5.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1755022273 with id 6.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1790522282 with id 5.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1785587280 with id 6.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1757271274 with id 7.
LogLearning: Display: LearningAgentsManager: Adding Agent BP_AICharacterB_C_UID_18C04D4052A589F201_1803042285 with id 7.
PIE: Server logged in
```

Fig.2 Instance starting

Fig.3 and Fig.4 Shows the initialization of the training process.

```
PIE: Server logged in
PIE: Play in editor total start time 0.764 seconds.
LogLearning: Display: TrainerB_0: Sending / Receiving initial policy...
LogLearning: Display: Training Process: {
LogLearning: Display: Training Process: "TaskName": "TrainerB_0",
LogLearning: Display: Training Process: "TrainerMethod": "PPO",
LogLearning: Display: Training Process: "TrainerType": "SharedMemory",
LogLearning: Display: Training Process: "TimeStamp": "2025-01-05_11-05-47",
LogLearning: Display: Training Process: "ExtraSitePackagesPath": "G:/UE_5.4/Engine/Binaries/Win64/",
LogLearning: Display: Training Process: "IntermediatePath": "G:/UE_5.4/Projects/LearnShooting/Intermediate/LearningAgents",
LogLearning: Display: Training Process: "PolicyGuid": "(22005B41-4C94-1496-02BE-1B8B69A7598A)",
LogLearning: Display: Training Process: "CriticGuid": "(632957F6-45C2-9C6E-A314-67B93FEC3440)",
LogLearning: Display: Training Process: "EncoderGuid": "(5839786C-4260-1253-72EE-67A40CF02318)",
LogLearning: Display: Training Process: "DecoderGuid": "(801F49E3-45FD-E989-5865-AFA8B74EC75C)",
LogLearning: Display: Training Process: "ControlsGuid": "(638FAFA8-4CDE-10E2-55D4-1DAF0C0F1B5)",
LogLearning: Display: Training Process: "EpisodeStartGuid": "(3E14E2F6-48F0-64F7-078B-0A8B34016FC3)",
LogLearning: Display: Training Process: "EpisodeLengthGuid": "(794684F4-430E-9560-446B-4A9E49A2A210)",
LogLearning: Display: Training Process: "EpisodeCompletionNodesGuid": "(C7AC9E0E-4A5E-94B7-8724-9E44100DE36B)",
LogLearning: Display: Training Process: "EpisodeFinalObservationsGuid": "(3EF18805-4461-2334-9E8B-00BFCAC60DE)",
LogLearning: Display: Training Process: "EpisodeFinalMemoryStatesGuid": "(44BFB6A4-4FD1-930F-3E55-BF8B78EF181)",
LogLearning: Display: Training Process: "ObservationsGuid": "(2633A136-4ACF-8CDB-2B6F-BE8EF470B800)",
LogLearning: Display: Training Process: "ActionsGuid": "(603788E3-49E3-0BC0-F8DA-529EF20E00C8)",
LogLearning: Display: Training Process: "MemoryStatesGuid": "(5383FB7E-4B56-54D4-B9C2-C6A4622C1E44)",
LogLearning: Display: Training Process: "RewardsGuid": "(DAAEC0C4-4BE7-6E84-6034-BD8681026AC3)",
```

Fig.3 Training Process initialization

```
Details OutputLog x
Search Log Filters v
PIE: Play in editor total start time 0.764 seconds.
LogLearning: Display: TrainerB_0: Sending / Receiving initial policy...
LogLearning: Display: Training Process: {
LogLearning: Display: Training Process: "TaskName": "TrainerB_0",
LogLearning: Display: Training Process: "TrainerMethod": "PPO",
LogLearning: Display: Training Process: "TrainerType": "SharedMemory",
LogLearning: Display: Training Process: "TimeStamp": "2025-01-05_11-05-47",
LogLearning: Display: Training Process: "ExtraSitePackagesPath": "G:/UE_5.4/Engine/Binaries/Win64/",
LogLearning: Display: Training Process: "IntermediatePath": "G:/UE_5.4/Projects/LearnShooting/Intermediate/LearningAgents",
LogLearning: Display: Training Process: "PolicyGuid": "(22005B41-4C94-1496-02BE-1B8B69A7598A)",
LogLearning: Display: Training Process: "CriticGuid": "(632957F6-45C2-9C6E-A314-67B93FEC3440)",
LogLearning: Display: Training Process: "EncoderGuid": "(5839786C-4260-1253-72EE-67A40CF02318)",
LogLearning: Display: Training Process: "DecoderGuid": "(801F49E3-45FD-E989-5865-AFA8B74EC75C)",
LogLearning: Display: Training Process: "ControlsGuid": "(638FAFA8-4CDE-10E2-55D4-1DAF0C0F1B5)",
LogLearning: Display: Training Process: "EpisodeStartGuid": "(3E14E2F6-48F0-64F7-078B-0A8B34016FC3)",
LogLearning: Display: Training Process: "EpisodeLengthGuid": "(794684F4-430E-9560-446B-4A9E49A2A210)",
LogLearning: Display: Training Process: "EpisodeCompletionNodesGuid": "(C7AC9E0E-4A5E-94B7-8724-9E44100DE36B)",
LogLearning: Display: Training Process: "EpisodeFinalObservationsGuid": "(3EF18805-4461-2334-9E8B-00BFCAC60DE)",
LogLearning: Display: Training Process: "EpisodeFinalMemoryStatesGuid": "(44BFB6A4-4FD1-930F-3E55-BF8B78EF181)",
LogLearning: Display: Training Process: "ObservationsGuid": "(2633A136-4ACF-8CDB-2B6F-BE8EF470B800)",
LogLearning: Display: Training Process: "ActionsGuid": "(603788E3-49E3-0BC0-F8DA-529EF20E00C8)",
LogLearning: Display: Training Process: "MemoryStatesGuid": "(5383FB7E-4B56-54D4-B9C2-C6A4622C1E44)",
LogLearning: Display: Training Process: "RewardsGuid": "(DAAEC0C4-4BE7-6E84-6034-BD8681026AC3)",
```

Fig.4 Initialization of policy, episode, critical, encoder, decoder and control guides

Fig.5, 6 & 7 showcases the starting iteration, second iteration and 480th iteration of the training process.

```
LogLearning: Display: Training Process: Receiving Policy...
LogLearning: Display: Training Process: Receiving Critic...
LogLearning: Display: Training Process: Receiving Encoder...
LogLearning: Display: Training Process: Receiving Decoder...
LogLearning: Display: Training Process: Creating Optimizer...
LogLearning: Display: Training Process: Creating PPO Policy...
LogLearning: Display: Training Process: Opening TensorBoard...
LogLearning: Display: Training Process: Begin Training...
LogLearning: Display: Training Process: Profile| Pull Experience      19297ms
LogLearning: Display: Training Process: Profile| PPO load tensors  1ms
LogLearning: Display: Training Process: Profile| PPO gae            35ms
LogLearning: Display: Training Process: Profile| PPO log prob       53ms
LogLearning: Display: Training Process: Profile| PPO create batches  44ms
LogLearning: Display: Training Process: Profile| PPO train          9087ms
LogLearning: Display: Training Process: Profile| Training          12268ms
LogLearning: Display: Training Process: Profile| Pushing Policy     4ms
LogLearning: Display: Training Process: Profile| Pushing Critic    1ms
LogLearning: Display: Training Process: Profile| Pushing Encoder    2ms
LogLearning: Display: Training Process: Profile| Pushing Decoder    1ms
LogLearning: Display: Training Process: Iter:      0 | Avg Rew: -0.31128 | Avg Rew Sum: -159.37500 |
```

Fig.5 First Iteration

```
LogDerivedDataCache: C:/Users/visal/AppData/Local/UnrealEngine/Common/DerivedDataCache: Maintenance fini
LogLearning: Display: LearningAgentsManager: Resetting Agents [0 1 2 3 4 5 6 7].
LogLearning: Display: Training Process: Profile| Logging           0ms
LogLearning: Display: Training Process: Profile| Pull Experience    28709ms
LogLearning: Display: Training Process: Profile| PPO load tensors  1ms
LogLearning: Display: Training Process: Profile| PPO gae            5ms
LogLearning: Display: Training Process: Profile| PPO log prob       2ms
LogLearning: Display: Training Process: Profile| PPO create batches  43ms
LogLearning: Display: Training Process: Profile| PPO train          6148ms
LogLearning: Display: Training Process: Profile| Training          6199ms
LogLearning: Display: Training Process: Profile| Pushing Policy     2ms
LogLearning: Display: Training Process: Profile| Pushing Critic    0ms
LogLearning: Display: Training Process: Profile| Pushing Encoder    1ms
LogLearning: Display: Training Process: Profile| Pushing Decoder    0ms
LogLearning: Display: Training Process: Iter:     32 | Avg Rew: -0.00342 | Avg Rew Sum: -1.75000 | Avg
```

Fig.6 Second Iteration

```
LogLearning: Display: Training Process: Iter:     416 | Avg Rew: 0.03418 | Avg Rew Sum: 17.5000
LogLearning: Display: LearningAgentsManager: Resetting Agents [0 1 2 3 4 5 6 7].
LogLearning: Display: Training Process: Profile| Logging           0ms
LogLearning: Display: Training Process: Profile| Pull Experience    20274ms
LogLearning: Display: Training Process: Profile| PPO load tensors  7ms
LogLearning: Display: Training Process: Profile| PPO gae            2ms
LogLearning: Display: Training Process: Profile| PPO log prob       46ms
LogLearning: Display: Training Process: Profile| PPO create batches  5803ms
LogLearning: Display: Training Process: Profile| PPO train          5859ms
LogLearning: Display: Training Process: Profile| Pushing Policy     2ms
LogLearning: Display: Training Process: Profile| Pushing Critic    2ms
LogLearning: Display: Training Process: Profile| Pushing Encoder    0ms
LogLearning: Display: Training Process: Profile| Pushing Decoder    1ms
LogLearning: Display: Training Process: Iter:     416 | Avg Rew: -0.15063 | Avg Rew Sum: -77.12
LogLearning: Display: Training Process: Profile| Logging           0ms
LogLearning: Display: Training Process: Profile| Pull Experience    13100ms
LogLearning: Display: Training Process: Profile| PPO load tensors  4ms
LogLearning: Display: Training Process: Profile| PPO gae            70ms
LogLearning: Display: Training Process: Profile| PPO log prob       34ms
LogLearning: Display: Training Process: Profile| PPO create batches  44ms
LogLearning: Display: Training Process: Profile| PPO train          8474ms
LogLearning: Display: Training Process: Profile| Training          8626ms
LogLearning: Display: Training Process: Profile| Pushing Policy     3ms
LogLearning: Display: Training Process: Profile| Pushing Critic    1ms
LogLearning: Display: Training Process: Profile| Pushing Encoder    2ms
LogLearning: Display: Training Process: Profile| Pushing Decoder    0ms
LogLearning: Display: Training Process: Iter:     448 | Avg Rew: 0.04004 | Avg Rew Sum: 20.5000
LogLearning: Display: LearningAgentsManager: Resetting Agents [0 1 2 3 4 5 6 7].
LogLearning: Display: LearningAgentsManager: Resetting Agents [0 1 2 3 4 5 6 7].
LogLearning: Display: Training Process: Profile| Logging           0ms
LogLearning: Display: Training Process: Profile| Pull Experience    15877ms
LogLearning: Display: Training Process: Profile| PPO load tensors  2ms
LogLearning: Display: Training Process: Profile| PPO gae            6ms
LogLearning: Display: Training Process: Profile| PPO log prob       2ms
LogLearning: Display: Training Process: Profile| PPO create batches  47ms
LogLearning: Display: Training Process: Profile| PPO train          5819ms
LogLearning: Display: Training Process: Profile| Training          5875ms
LogLearning: Display: Training Process: Profile| Pushing Policy     2ms
LogLearning: Display: Training Process: Profile| Pushing Critic    1ms
LogLearning: Display: Training Process: Profile| Pushing Encoder    0ms
LogLearning: Display: Training Process: Profile| Pushing Decoder    1ms
LogLearning: Display: Training Process: Iter:     448 | Avg Rew: -0.64615 | Avg Rew Sum: -313.6
LogLearning: Display: Training Process: Profile| Logging           0ms
LogLearning: Display: Training Process: Profile| Pull Experience    13510ms
LogLearning: Display: Training Process: Profile| PPO load tensors  4ms
LogLearning: Display: Training Process: Profile| PPO gae            120ms
LogLearning: Display: Training Process: Profile| PPO log prob       35ms
LogLearning: Display: Training Process: Profile| PPO create batches  43ms
LogLearning: Display: Training Process: Profile| PPO train          8738ms
LogLearning: Display: Training Process: Profile| Training          8960ms
LogLearning: Display: Training Process: Profile| Pushing Policy     2ms
LogLearning: Display: Training Process: Profile| Pushing Critic    1ms
LogLearning: Display: Training Process: Profile| Pushing Encoder    2ms
LogLearning: Display: Training Process: Profile| Pushing Decoder    0ms
LogLearning: Display: Training Process: Iter:     480 | Avg Rew: 0.11251 | Avg Rew Sum: 37.7500
LogLearning: Display: LearningAgentsManager: Resetting Agents [0 1 2 3 4 5 6 7].
LogLearning: Display: LearningAgentsManager: Resetting Agents [0 1 2 3 4 5 6 7].
LogLearning: Display: Training Process: Profile| Logging           0ms
LogLearning: Display: Training Process: Profile| Pull Experience    16590ms
LogLearning: Display: Training Process: Profile| PPO load tensors  1ms
LogLearning: Display: Training Process: Profile| PPO gae            8ms
LogLearning: Display: Training Process: Profile| PPO log prob       3ms
LogLearning: Display: Training Process: Profile| PPO create batches  4ms
LogLearning: Display: Training Process: Profile| PPO train          5716ms
LogLearning: Display: Training Process: Profile| Training          5773ms
LogLearning: Display: Training Process: Profile| Pushing Policy     3ms
LogLearning: Display: Training Process: Profile| Pushing Critic    0ms
LogLearning: Display: Training Process: Profile| Pushing Encoder    1ms
LogLearning: Display: Training Process: Profile| Pushing Decoder    0ms
LogLearning: Display: Training Process: Iter:     480 | Avg Rew: 0.30396 | Avg Rew Sum: 155.625
```

Fig.7 480th iteration

Fig. 8 (a, b, c & d) showcases the in-game scoreboard for both the teams. This clearly showcases how both teams have won the games throughout the training process instead of only one team getting the upper hand in the score.

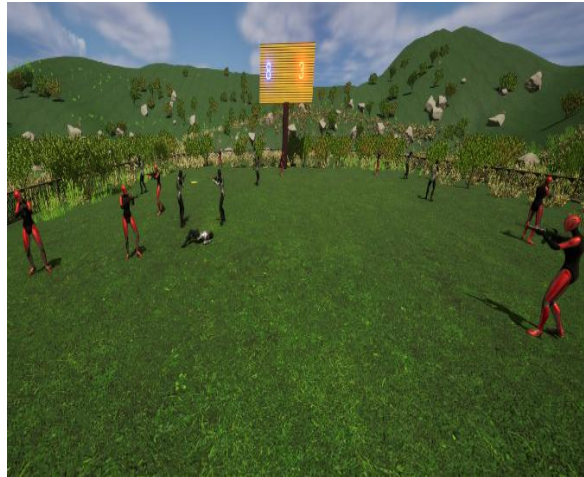


Fig.8 (a) Team A and Team B with scoreboard



Fig. 8 (b) Team B with more points



Fig. 8 (c) Team A with more points



Fig. 8 (d) Team A winning

Fig.9 Asset Validation after training process

```
LogSavePackage: Moving output files for package: /Game/Learn/Red/DA_NeuralNetworkPoli
LogSavePackage: Moving : ../../../../../../UE5.4/Projects/LearnShooting/Saved/DA_NeuralNetwo
LogFileHelpers: InternalPromptForCheckoutAndSave took 134,384 ms
LogContentValidation: Display: Starting to validate 8 assets
LogContentValidation: Enabled validators:
LogContentValidation: /Script/DataValidation.DirtyFilesChangelistValidator
LogContentValidation: /Script/DataValidation.EditorValidator_Localization
LogContentValidation: /Script/DataValidation.WorldPartitionChangelistValidator
LogContentValidation: Display: Validating asset /Game/Learn/Blue/DA_NeuralNetworkEnco
AssetCheck: /Game/Learn/Blue/DA_NeuralNetworkEncoderB Validating asset
LogContentValidation: Display: Validating asset /Game/Learn/Blue/DA_NeuralNetworkDeco
AssetCheck: /Game/Learn/Blue/DA_NeuralNetworkDecoderB Validating asset
LogContentValidation: Display: Validating asset /Game/Learn/Blue/DA_NeuralNetworkCrit
AssetCheck: /Game/Learn/Blue/DA_NeuralNetworkCriticB Validating asset
LogContentValidation: Display: Validating asset /Game/Learn/Blue/DA_NeuralNetworkPoli
AssetCheck: /Game/Learn/Blue/DA_NeuralNetworkPolicyB Validating asset
LogContentValidation: Display: Validating asset /Game/Learn/Red/DA_NeuralNetworkEncod
AssetCheck: /Game/Learn/Red/DA_NeuralNetworkEncoderR Validating asset
LogContentValidation: Display: Validating asset /Game/Learn/Red/DA_NeuralNetworkDecod
AssetCheck: /Game/Learn/Red/DA_NeuralNetworkDecoderR Validating asset
LogContentValidation: Display: Validating asset /Game/Learn/Red/DA_NeuralNetworkCritic
AssetCheck: /Game/Learn/Red/DA_NeuralNetworkCriticR Validating asset
LogContentValidation: Display: Validating asset /Game/Learn/Red/DA_NeuralNetworkPolic
AssetCheck: /Game/Learn/Red/DA_NeuralNetworkPolicyR Validating asset
```

Fig.9 depicts the asset validation at the end of the training process.

Discussion

Key Benefits of PPO Architecture

The RL Agent Manager may periodically instruct the system to reset episodes, change the conditions for domain randomization. It collects global stats and manages training schedules by stopping once performance targets are met.

- i. **Scalability:** Multiple Unreal Engine instances can be distributed across several machines, each running a subset of agents which will effectively parallelize collected experience.
- ii. **Stable Training:** PPO's clipped updates in a centralized server will prevent catastrophic forgetting thus ensuring more robust learning in a competitive multi-agent setting.
- iii. **Modularity:** The separation between Unreal Engine environment and the Parameter Server learning provides flexible substitution of RL algorithm and parameter tuning scripts.
- iv. **Sample Efficiency:** Collected trajectories can be reused multiple times within a single update epoch which makes PPO effectively extract information from each sample thus resulting in accelerating learning in high-dimensional multi-agent environments. Each agent independently uses the learned policy which results in fostering scalability and robustness to partial observability or network latency.

Limitations

When running multiple Unreal Engine instances, each simulating multiple PPO agents, is a huge resource intensive process. Choices such as the clipping parameter ϵ , learning rate, and reward weights can substantially affect training outcomes due to hyper parameter sensitivity. Agents cannot always have a comprehensive view of the battlefield due to lack of deeper research into recurrence (RNN-based PPO) or attention mechanisms.

V. CONCLUSION

This paper about Adaptive multiple AI gave a hands-on approach on how Proximal Policy Optimization (PPO) can powerfully drive a multi-agent combat system in Unreal Engine, with two opposing teams striving to eliminate each other. By carefully designing state representations, action spaces, and reward structures, agents can learn coordinated behaviors and complex tactics. Experimental results show that PPO-trained agents offer a more engaging and dynamic challenge compared to traditional rule-based hard coded NPCs, indicating the potential for broader applications of RL in team-based adversarial game environments. Although resource demands and hyperparameter tuning remain substantial challenges, ongoing work in hierarchical methods and advanced sensor models may further elevate the sophistication of AI combat behavior. As the field evolves, RL-driven multi-agent systems will become more prevalent in both research and commercial gaming which offers players rich, adaptive experiences that clearly mimic real human behavior.

Future Directions

Splitting the combat tasks into higher-level strategies and lower-level combat skills will boost learning efficiency. New approaches that allow policies to generalize quickly to new map layouts or novel enemy strategies. Introducing new objectives such as resource collection or control point captures to blend cooperative and competitive dynamics. Using Leveraging advanced sensor models to identify and track enemies over long distances or through obstructed lines of sight.

References

- [1] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, O. Klimov (2017). “Proximal Policy Optimization Algorithms.”
- [2] David Silver, Aja Huang., (2016). “Mastering the game of Go with deep neural networks and tree search.”
- [3] O. Vinyals, T. Ewalds. (2019). “Grandmaster level in StarCraft II using multi-agent reinforcement learning.”
- [4] Remi Niel “Hierarchical Reinforcement Learning for Real-Time Strategy Games”
- [5] E. Lockhart and Y. K Kwok, (2021). “Resource Management Using Proximal Policy Optimization in Tower Defense Games.”
- [6] L. Zhang, A. A. Aziz, and M. G. Rashed, (2022). “Multi-Agent Deep Reinforcement Learning for Dynamic Team Tactics in MOBA Games.”
- [7] K. Kim, S. Im, J. Park, (2022) “Applying Resinforcement Learning in FPS Environments for Realistic Non-Player Characters.”
- [8] K. Shao, Z. Tang (2019). “A Survey of Deep Reinforcement Learning in Video Games.”
- [9] D. Kreffer, (2022). “Reinforcement learning for non-player characters in the video game industry.” Epic Games. Unreal Engine 5