

# Capture The Flag (CTF) using Multi Agent RL Algorithms

Shirish Shekhar Jha

Department of Electrical Engineering and Computer Science  
IISER Bhopal

November 25, 2024

# Overview

- 1 Problem Statement
- 2 Algorithms
- 3 Markov Decision Process (MDP)
- 4 Results
- 5 Results and Comparisons

# Problem Statement

- Capture the Flag environment has two teams with two agents in each team.
- Every team has the objective of capturing the opponent's flag, but at the same time defend its own.
- Defending the flag activates when an agent enters a visual depth of 3 near the opponent's flag.
- Obstacles, and flags positions were static, and two agents could occupy same cell.

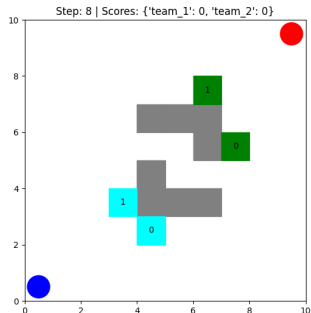


Figure: CTF Environment

# Reinforcement Learning

Below are the algorithms used in the study:-

- Independent Q-Learning
- Multi Agent Proximal Policy Optimization (**Preferred Approach**)

# MDP: Definition and Components (Part 1)

A Markov Decision Process (MDP) is defined as the tuple  $(S, A, d_R, d_0, \gamma)$ , where:

- $S$  denotes the state space. For our environment:

$$S \in \mathbb{R}^{10 \times 10 \times 4} \setminus S'$$

where  $S'$  represents the spaces occupied by obstacles or flag positions.

- $A$  defines the action space:

$$A = \{\text{Up } (0, 1), \text{ Down } (0, -1), \text{ Left } (-1, 0), \text{ Right } (1, 0), \text{ Stay } (0, 0)\}$$

- $d_R$  represents the reward distribution. For our problem:
  - +100: For capturing the opponent's flag.
  - +25: For successfully defending the flag.
  - -25: For getting caught while intruding.
  - -2: For staying in the same position.

# MDP: Definition and Components (Part 2)

Continuing the tuple  $(S, A, d_R, d_0, \gamma)$ :

- $d_0$  (Initial State): Defines the starting positions of agents:

- Team 1:

$$x \sim \text{Uniform}(0, \frac{\text{grid\_size}[n_R]}{2}), y \sim \text{Uniform}(0, \frac{\text{grid\_size}[n_C]}{2}), (x, y) \notin S'$$

- Team 2:

$$x \sim \text{Uniform}(\frac{\text{grid\_size}[n_R]}{2}, n_R), y \sim \text{Uniform}(\frac{\text{grid\_size}[n_C]}{2}, n_C),$$

$$(x, y) \notin S', n_R, \& n_C$$

are number of rows and columns in grid.

- $\gamma$  (Discount Factor): Balances immediate and future rewards:

$$\gamma = 0.99$$

# Challenges in the Problem

- Effectively shaping reward for exploration, defending, and capturing
- Achieving team objectives, when to start exploring to capture, when to defend the own territory.
- Delayed rewards for defending made it challenging for agents to learn.

# Why Choose MAPPO?

- **Centralized Training with Decentralized Execution:** Facilitates effective coordination among agents.
- **Proven Performance:** Achieves competitive or superior results in cooperative multi-agent scenarios.
- **Stable Learning Dynamics:** On-policy nature ensures stability in complex interactions.



# Alignment with Problem Requirements

- **Coordination:** Enables agents to learn joint policies for balanced offensive and defensive strategies.
- **Stability:** Ensures stable learning in environments with complex agent interactions.

# Model Parameters

- **Policy Network:**

- **Input Layer:** 100 neurons (corresponding to the flattened observation space).
- **Hidden Layers:** Two fully connected layers with 128 neurons each, activated by ReLU functions.
- **Output Layer:** 5 neurons representing the action logits.

- **Optimizer:** Adam optimizer with a learning rate of 0.0003.

- **Policy Loss:** Clipped surrogate objective to ensure stability during training:

$$L_{\text{policy}} = -\mathbb{E} [\min (r_t(\theta)A, \text{clip} (r_t(\theta), 1 - \epsilon, 1 + \epsilon) A)]$$

where  $r_t(\theta) = \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}$  is the probability ratio.

# Defining Metrics for Evaluation

- **High Win Rate:** The algorithm with a higher win rate indicates better team performance.
- **High Average Scores:** Reflects consistent ability to achieve objectives.
- **Low Draw Rate:** Indicates decisive outcomes, less stalemates.
- **Performance Stability:** Variance in scores across episodes to measure consistency.

## Why These Metrics?

- Assess overall dominance and effectiveness of each algorithm.
- Lower variance in scores suggests consistent performance.

# Training Metrics: Rewards (IQL)

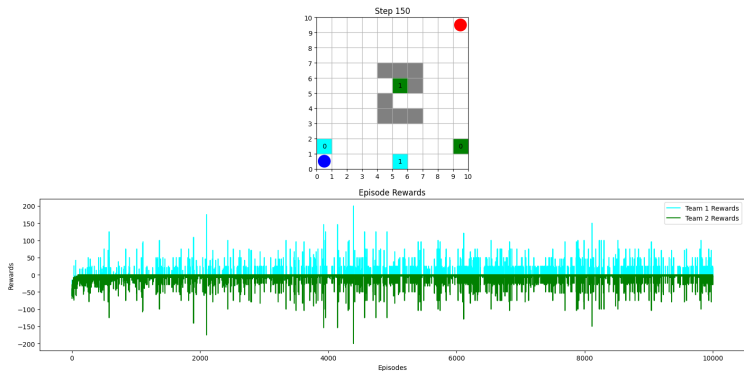


Figure: Loss progression during training (IQL)

# Training Metrics: Loss and Rewards (MAPPO)

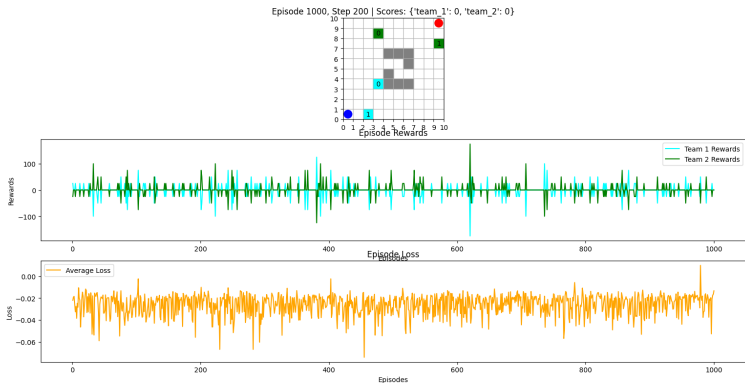
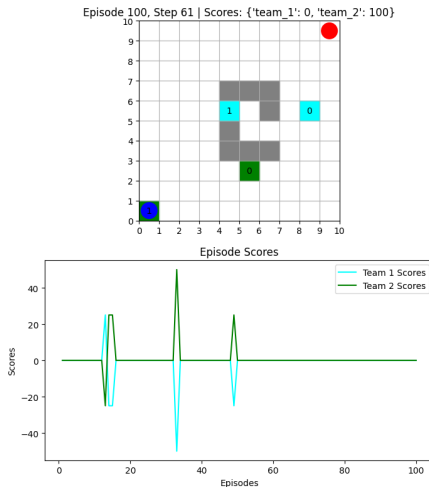


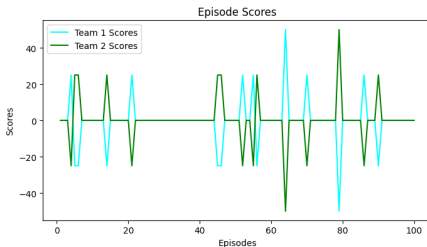
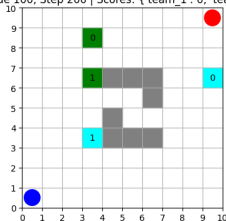
Figure: Loss progression during training (MAPPO)

# Testing Metrics: Rewards (IQL)



# Testing Metrics: Loss and Rewards (MAPPO)

Episode 100, Step 200 | Scores: {'team\_1': 0, 'team\_2': 0}



# Results and Observations

## Observations:

- **Win and Draw Rates:**

- MAPPO: Team 1 (53.13%), Team 2 (46.88%), Draws (0%).
- IQL: Team 1 (57.14%), Team 2 (42.86%), Draws (0%).

- **Average Scores:**

- MAPPO: Team 1 (-0.78125), Team 2 (0.78125).
- IQL: Team 1 (0.0), Team 2 (0.0).

- **Average Score Difference:**

- MAPPO: -1.5625, indicating stronger dynamics between teams.
- IQL: 0.0, demonstrating balanced team performance.

- MAPPO reflects greater variability in team performance due to centralized policy training.



# Results and Observations contd..

## Visualized Results:

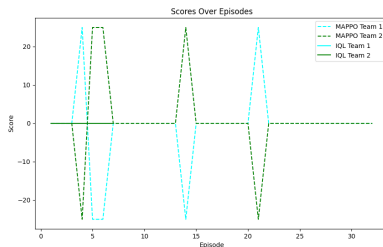
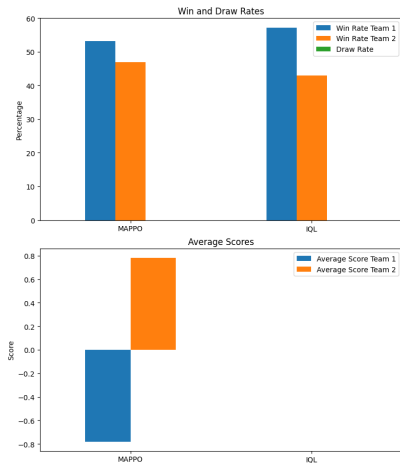


Figure: (Left) Win and Draw Rates; (Right) Scores Over Episodes.

# Environment Comparison

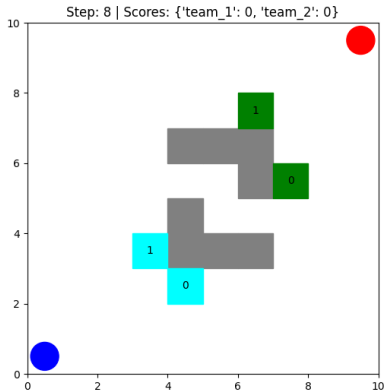


Figure: CTF Environment (Original)

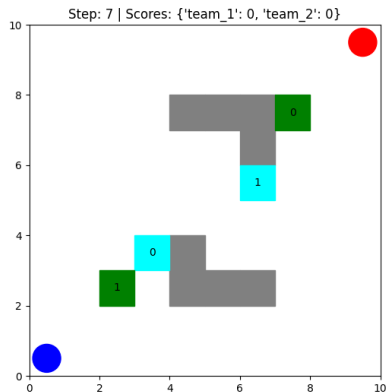


Figure: CTF Environment (Changed Obstacles)

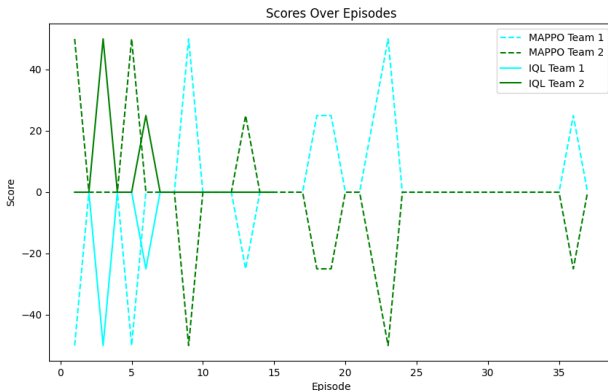
# Performance Metrics: Changed Environment

- **Win Rates:** MAPPO shows competitive balance (48.65% Team 1, 51.35% Team 2), while IQL favors Team 1 (53.33%).
- **Average Scores:** MAPPO results in balanced scores (2.03 for Team 1, -2.03 for Team 2), whereas IQL exhibits a significant score disparity (-5.00 for Team 1, 5.00 for Team 2).
- **Inferences:**
  - MAPPO's centralized coordination leads to balanced gameplay.
  - IQL's independent strategies result in unbalanced performance across teams.

# Performance Metrics Over Episodes

- **Score Evolution:**

- MAPPO exhibits stable dynamics with alternating scores over episodes, showing competitive engagement.
- IQL has steeper score fluctuations, indicating independent decisions often fail to adjust dynamically.



# Score Dynamics Over Episodes

## Visualized Results:

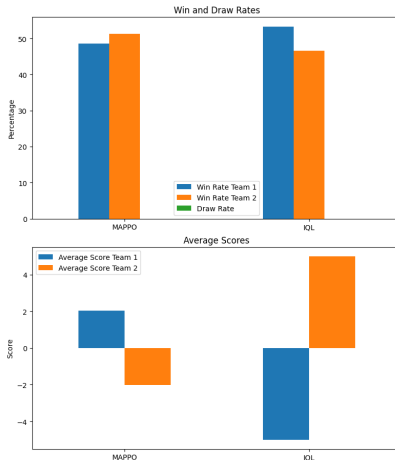
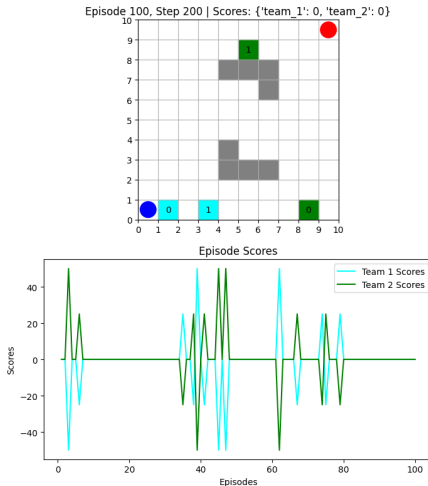


Figure: Win and Draw Rates

# Testing Metrics: Rewards (IQL) in Changed Environment



# Testing Metrics: Loss and Rewards (MAPPO) in Chaged Environment

