

# CSE 676 Deep Learning

## FINAL PROJECT

On

# FOREST FIRE MITIGATION

Using MARL (Multi-Agent  
Reinforcement Learning)

 University at Buffalo  
Information Technology

By:  
Group - 10  
Tanwin Chowdary (tanwinch)  
Kumar Raja Chidella (kumarraj)  
Veekshith Kajeepeta (vkajeepe)



# Project Description

- Objective: Develop an efficient solution for managing forest fires by utilizing a reinforcement learning based approach.
- Create a grid-based simulation environment that accurately models the dynamics of forest fires and allows for the simulation of various scenarios.
- Use Multiple agents, which can be firefighters and drones, will be deployed to move around the grid and take appropriate actions.
- By training RL algorithms, we aim to enable these agents to learn optimal firefighting strategies tailored to different situations
- Ultimately, minimize the damage caused by forest fires, safeguard lives, and enhance overall firefighting efficiency.



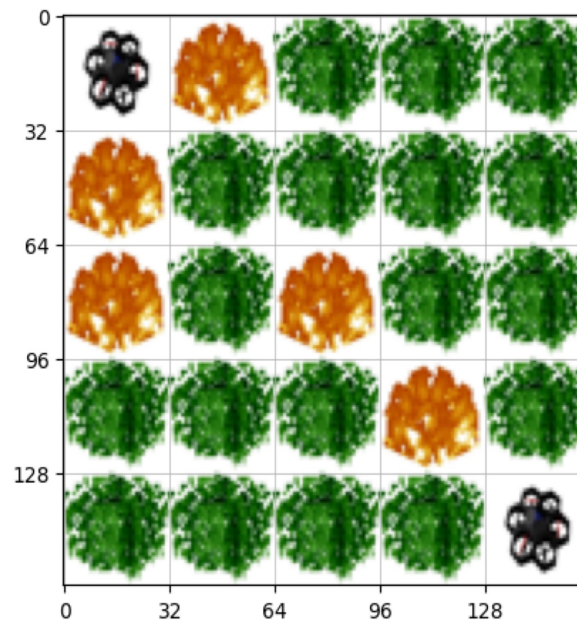
# Background

- In today's world, there is a worldwide problem of Wild Forest Fires and how to deal with it
- Need for mitigation of problem
- Replicate similar conditions in a test environment
- Integrate Reinforcement Learning Techniques to solve environment.
- Previous works show implementation of Forest Fire detection using CNN algorithms and Image classification.



# About the Environment

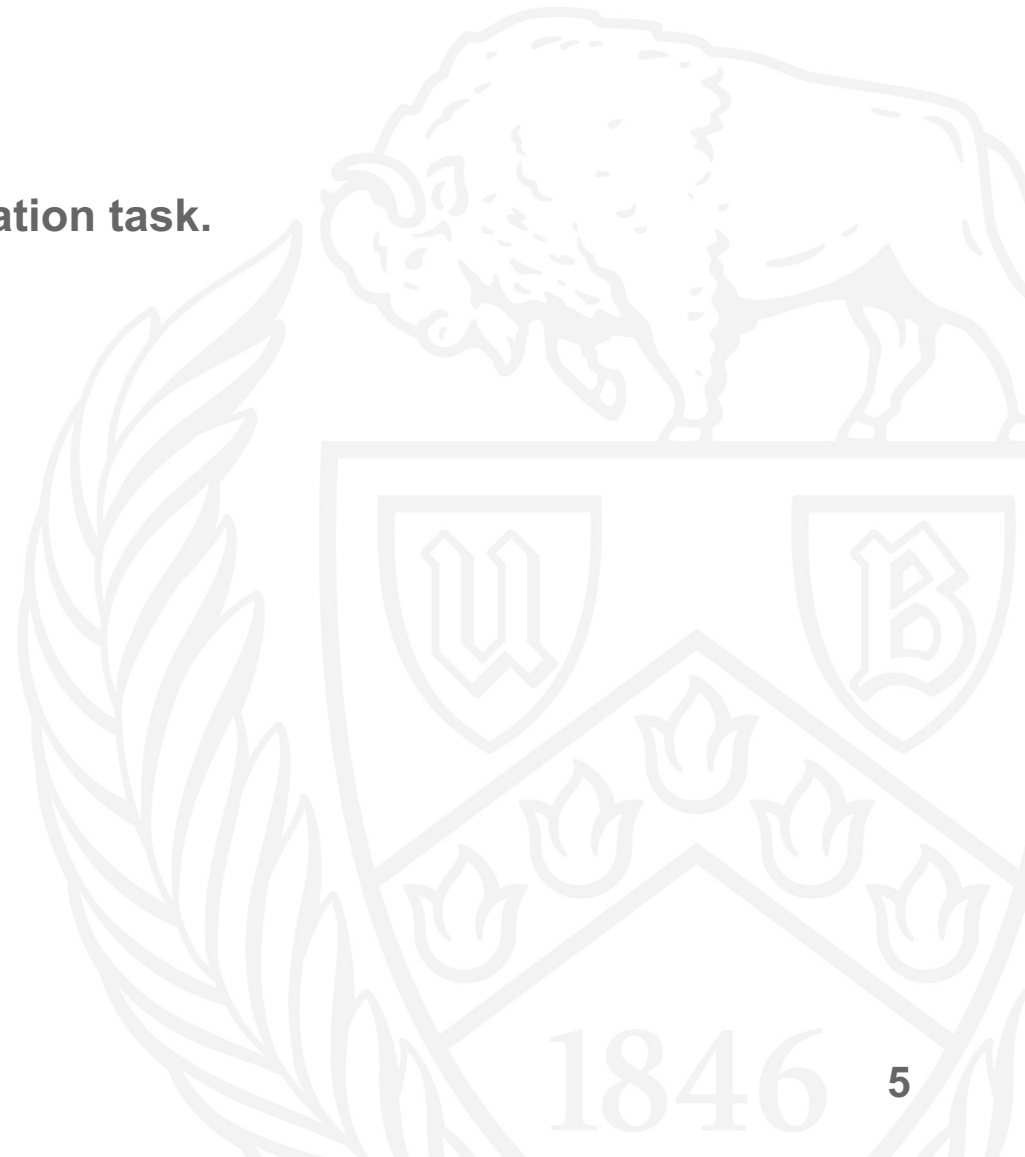
- Existing environments such as pendulum and CartPole from gym.
- Observation Space/Grid  
Dimensions:  $n \times n$
- Actions: Move Up/Down/Left/Right
- Rewards:
  - a) -0.1 for visiting a cell with no fire.
  - b) +2.0 for visiting a cell with fire and dousing it.
  - c) +5.0 for dousing all cells with fire.



- Has Stochastic and Deterministic variability
- Stochastic Environment
  - a) Random Position Fires.
  - b) Fires Spread with Random Probability.
- Deterministic Environment
  - a) Fire at Static Positions.
  - b) No spread of Fire.
- - Q-table with grid position where particular agent is in state.
- - Q-table with values of all grid positions in state.

# Methods

- The following methods were used to enhance forest fire mitigation task.
- Proximal Policy Optimisation (PPO)
- Soft Actor Critic (SAC)
- Deep Deterministic Policy Gradient (DDPG)
- Deep Q-Networks (DQN)
- Dueling Deep Q-Networks (DDQN)





## DQN

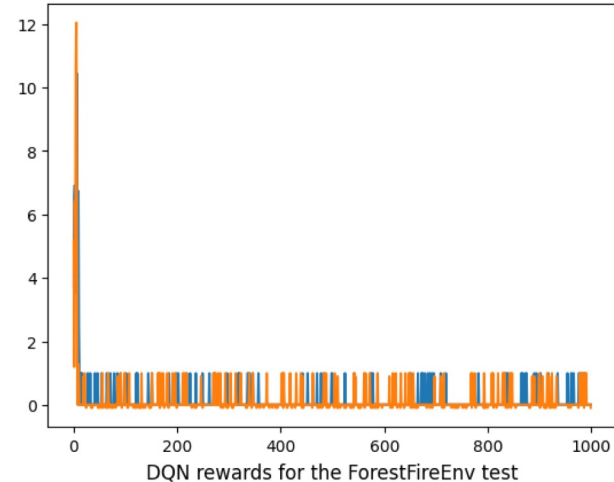
- DQNAgent is a class that represents an agent in the DQN algorithm.
- The agent has a neural network model that predicts the Q-values of each possible action for a given state.
- The agent uses an experience replay buffer to store experiences (state, action, reward, next\_state, done) during training.
- The agent chooses an action using an epsilon-greedy policy, where it selects a random action with probability  $\epsilon$  and selects the action with the highest Q-value with probability  $(1-\epsilon)$ .
- The agent trains its neural network by sampling a random batch of experiences from the replay buffer and updating the model's parameters using the Bellman equation.
- The agent updates the target network(which has the same architecture as the model but with frozen parameters, to estimate the Q-values used in the Bellman equation.) every  $\text{update\_frequency}=100$  steps.
- The code trains two DQN agents on the ForestFireEnv environment

### DQN

```
[ ] 1 import torch
    2 import torch.nn as nn
    3 import torch.optim as optim
    4 import random
    5 from collections import deque
    6
    7 class DQNAgent:
    8     def __init__(self, state_size, action_size, device, lr=0.001, gamma=0.99, buffer_size=10000, batch_size=64):
    9         self.state_size = state_size
   10         self.action_size = action_size
   11         self.device = device
   12         self.lr = lr
   13         self.gamma = gamma
   14         self.batch_size = batch_size
   15
   16         self.memory = deque(maxlen=buffer_size)
   17         self.model = self._build_model().to(device)
   18         self.target_model = self._build_model().to(device)
   19         self.optimizer = optim.Adam(self.model.parameters(), lr=self.lr)
   20         self.update_target_model()
   21
   22     def _build_model(self):
   23         model = nn.Sequential(
   24             nn.Linear(self.state_size, 64),
   25             nn.ReLU(),
   26             nn.Linear(64, 64),
   27             nn.ReLU(),
   28             nn.Linear(64, self.action_size)
   29         )
   30         return model
   31
   32     def remember(self, state, action, reward, next_state, done):
   33         self.memory.append((state, action, reward, next_state, done))
   34
   35     def choose_action(self, state, epsilon):
   36         if random.random() < epsilon:
   37             return random.randrange(self.action_size)
   38         else:
   39             state = torch.tensor(state, dtype=torch.float32, device=self.device).unsqueeze(0)
   40             return self.model(state).argmax().item()
   41
   42     def update_target_model(self):
   43         self.target_model.load_state_dict(self.model.state_dict())
   44
   45     def train(self):
   46         if len(self.memory) < self.batch_size:
   47             return
   48
   49         minibatch = random.sample(self.memory, self.batch_size)
   50
   51         states, actions, rewards, next_states, dones = zip(*minibatch)
   52         states = torch.tensor(states, dtype=torch.float32, device=self.device)
   53         actions = torch.tensor(actions, dtype=torch.long, device=self.device).unsqueeze(1)
   54         rewards = torch.tensor(rewards, dtype=torch.float32, device=self.device)
   55         next_states = torch.tensor(next_states, dtype=torch.float32, device=self.device)
   56         dones = torch.tensor(dones, dtype=torch.float32, device=self.device)
   57
   58         current_q_values = self.model(states).gather(1, actions)
   59         next_q_values = self.target_model(next_states).max(1)[0].detach()
   60         target_q_values = rewards + self.gamma * next_q_values * (1 - dones)
   61
   62         loss = nn.MSELoss()(current_q_values, target_q_values.unsqueeze(1))
   63         self.optimizer.zero_grad()
   64         loss.backward()
   65         self.optimizer.step()
```

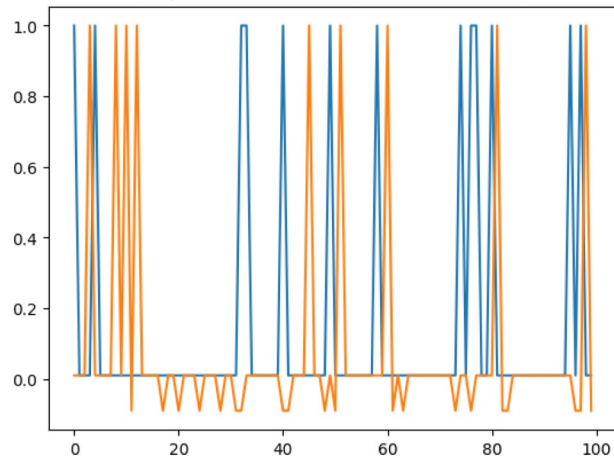
# DQN stochastic v/s deterministic

DQN rewards for the ForestFireEnv



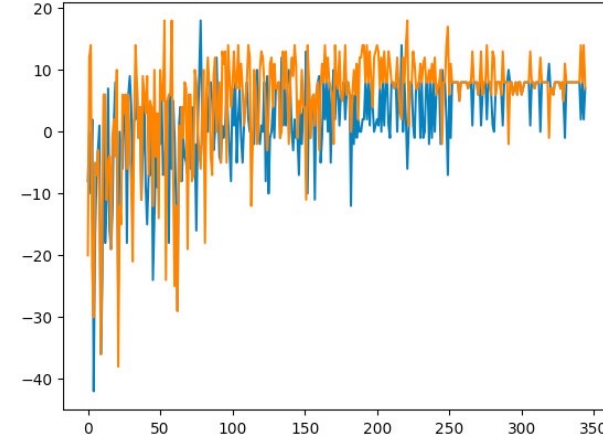
Due to the volatility of the environment defined the results of the DQN implemented are no as expected with max timesteps of 1000 on a 6x6 grid env the results are  $[-0.1, -0.1]$  for the agents.

DQN rewards for the ForestFireEnv test



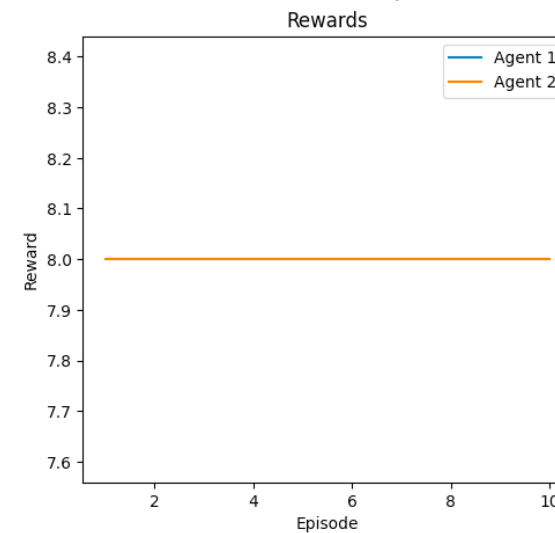
And because of the training handling the volatile env the test results are fluctuating for the stochastic env for the 100 episodes.

DQN rewards for the ForestFireEnv



Due to the stability of the environment defined the results of the DQN implemented are quite impressive with max timesteps of 100 on a 4x4 grid env the results are stabilizing at  $[8, 8]$  for the agents.

DQN Test Results:



And because of the training handling the stable env the test results are constant for the static env for the 10 episodes. Results are achieved in 4 timesteps.

# DDQN

- It is an improvement over DQN, which reduces overestimation of Q-values.
- It involves using two separate neural networks to estimate Q-values: the online model and the target model.
- The online model is used to select actions, while the target model is used to evaluate the actions.
- In DDQN, the target Q-values are computed using the target model, but the action selection is done using the online model.
- The code for DDQNAgent is similar to that of DQNAgent, but with modifications in the train() method to compute the target Q-values using the target model and the online model.
- The DDQN code presented here trains the agents using the ForestFireEnv.
- The training loop is similar to that of DQN, but with the use of two separate models to estimate Q-values.

## DDQN

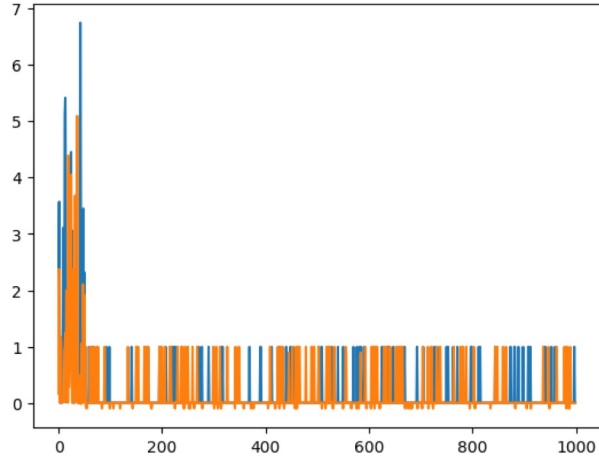
```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import random
5 from collections import deque
6
7 class DDQNAgent:
8     def __init__(self, state_size, action_size, device, lr=0.001, gamma=0.99, buffer_size=10000, batch_size=64):
9         self.state_size = state_size
10        self.action_size = action_size
11        self.device = device
12        self.lr = lr
13        self.gamma = gamma
14        self.batch_size = batch_size
15
16        self.memory = deque(maxlen=buffer_size)
17        self.model = self._build_model().to(device)
18        self.target_model = self._build_model().to(device)
19        self.optimizer = optim.Adam(self.model.parameters(), lr=self.lr)
20        self.update_target_model()
21
22    def _build_model(self):
23        model = nn.Sequential(
24            nn.Linear(self.state_size, 64),
25            nn.ReLU(),
26            nn.Linear(64, 64),
27            nn.ReLU(),
28            nn.Linear(64, self.action_size)
29        )
30        return model
31
32    def remember(self, state, action, reward, next_state, done):
33        self.memory.append((state, action, reward, next_state, done))
34
35    def choose_action(self, state, epsilon):
36        if random.random() < epsilon:
37            return random.randrange(self.action_size)
38        else:
39            state = torch.tensor(state, dtype=torch.float32, device=self.device).unsqueeze(0)
40            return self.model(state).argmax().item()
41
42    def update_target_model(self):
43        self.target_model.load_state_dict(self.model.state_dict())
44
45    def train(self):
46        if len(self.memory) < self.batch_size:
47            return
48
49        minibatch = random.sample(self.memory, self.batch_size)
50
51        states, actions, rewards, next_states, dones = zip(*minibatch)
52        states = torch.tensor(states, dtype=torch.float32, device=self.device)
53        actions = torch.tensor(actions, dtype=torch.long, device=self.device).unsqueeze(1)
54        rewards = torch.tensor(rewards, dtype=torch.float32, device=self.device)
55        next_states = torch.tensor(next_states, dtype=torch.float32, device=self.device)
56        dones = torch.tensor(dones, dtype=torch.float32, device=self.device)
57
58        current_q_values = self.model(states).gather(1, actions)
59
60        # Use online model to select actions
61        online_next_actions = self.model(next_states).argmax(1).unsqueeze(1)
62        # Use target model to evaluate actions
63        next_q_values = self.target_model(next_states).gather(1, online_next_actions).squeeze(1).detach()
64
65        target_q_values = rewards + self.gamma * next_q_values * (1 - dones)
66
67        loss = nn.MSELoss()(current_q_values, target_q_values.unsqueeze(1))
68        self.optimizer.zero_grad()
69        loss.backward()
70        self.optimizer.step()
  
```



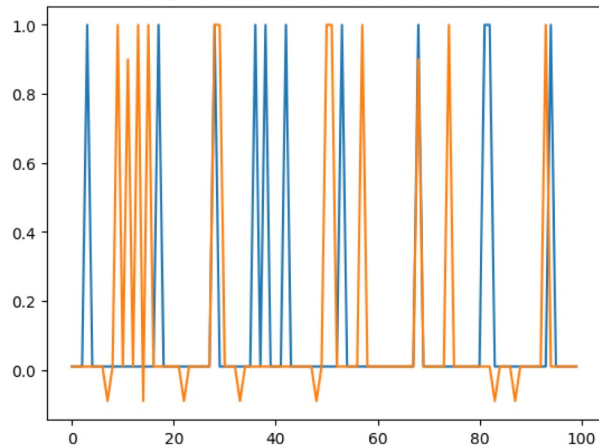
# DDQN stochastic v/s deterministic

DDQN rewards for the ForestFireEnv



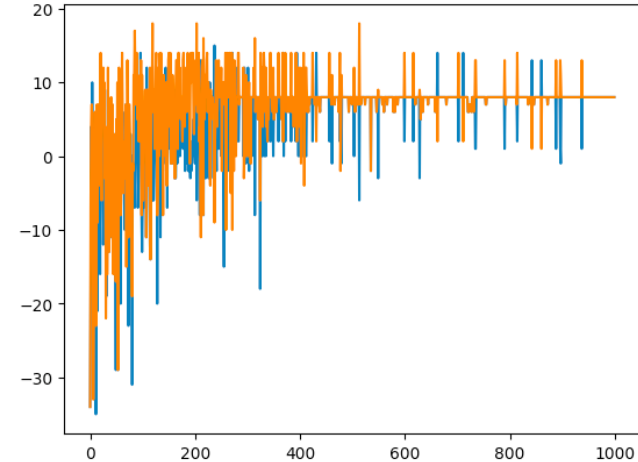
Due to the volatility of the environment defined the results of the DDQN implemented are no as expected with max timesteps of 1000 on a 6x6 grid env the results are [0.1, 1] for the agents.

DDQN rewards for the ForestFireEnv test



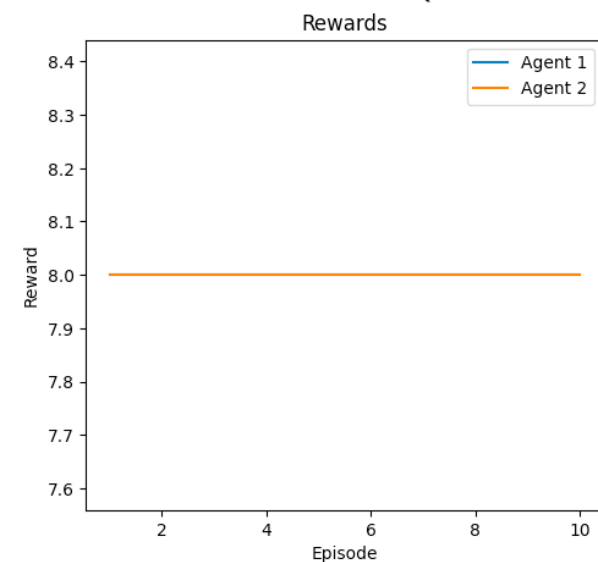
And because of the training handling the volatile env the test results are fluctuating for the stochastic env for the 100 episodes.

DDQN rewards for the ForestFireEnv



Due to the stability of the environment defined the results of the DDQN implemented are quite impressive with max timesteps of 100 on a 4x4 grid env the results are stabilizing at [8, 8] for the agents.

DDQN Test Results:

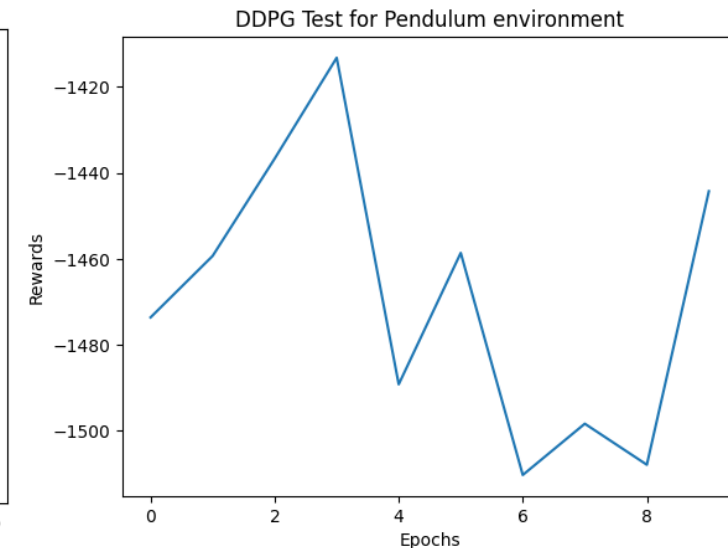
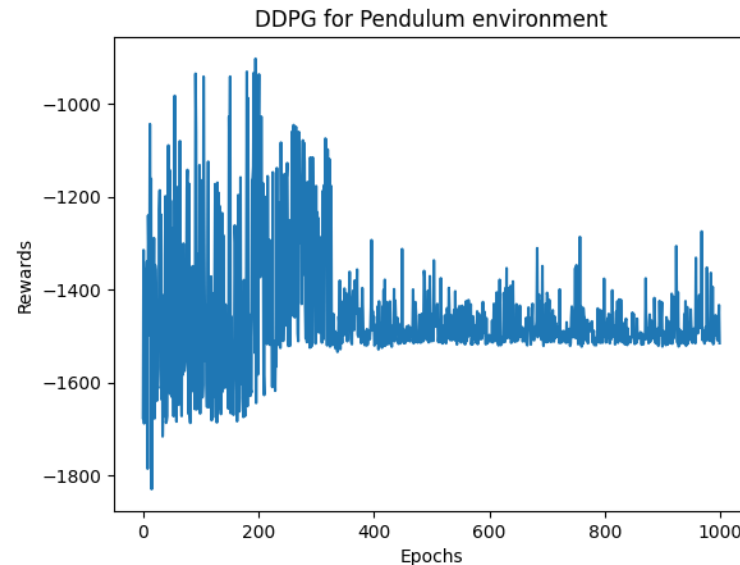
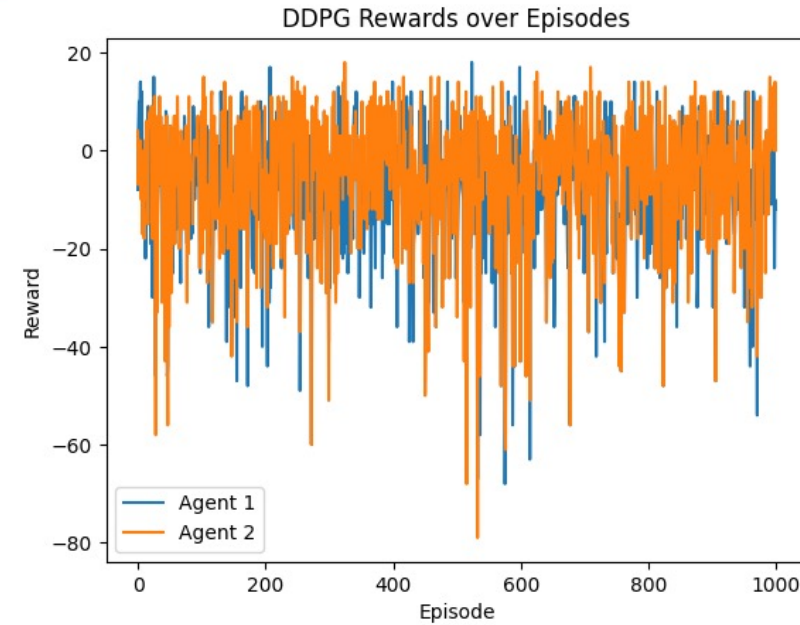


And because of the training handling the stable env the test results are constant for the static env for the 10 episodes. Results are achieved in 4 timesteps.

# Results

## DDPG:

- The above graph represents the DDPG performance on our stochastic environment, where fire starts and spreads randomly on the grid.
- The below following graphs show us the implementation of the DDPG algorithm on the pendulum gym environment. As we can see that after a certain number of epochs, the model has stabilized with some satisfactory rewards.

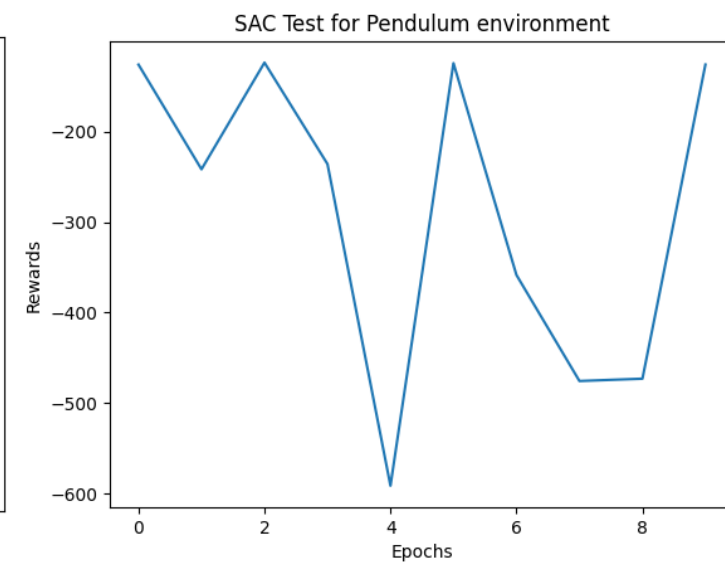
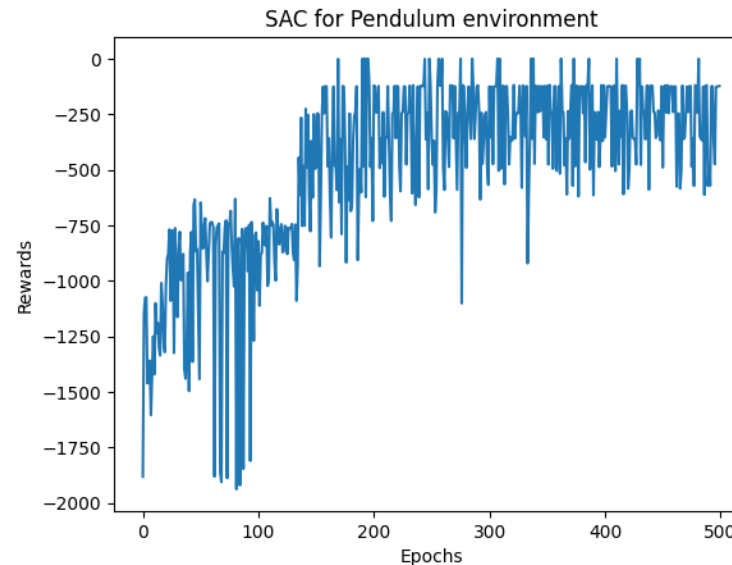
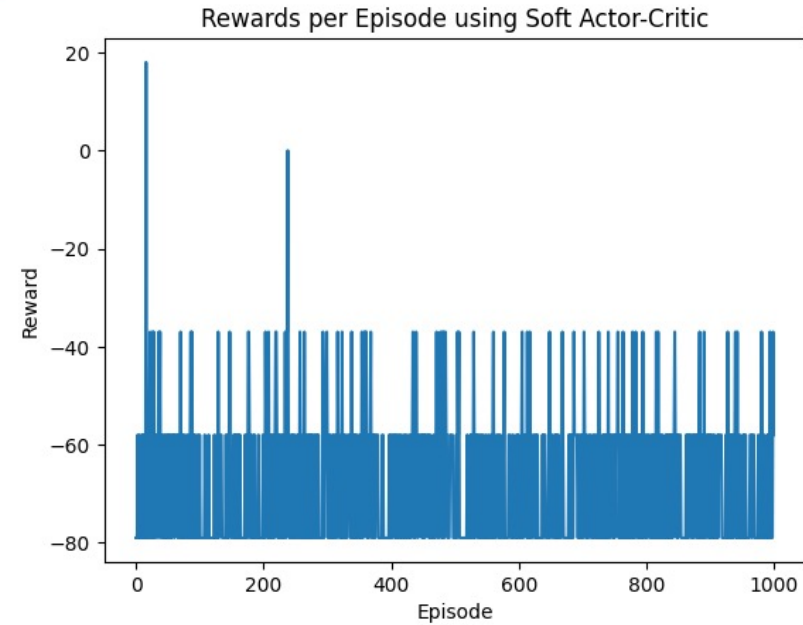


# Results

## SAC:

The above graph represents the SAC algorithm on the stochastic ForestFireEnv, SAC performance has been pretty underwhelming in this environment compared to the pendulum gym environment.

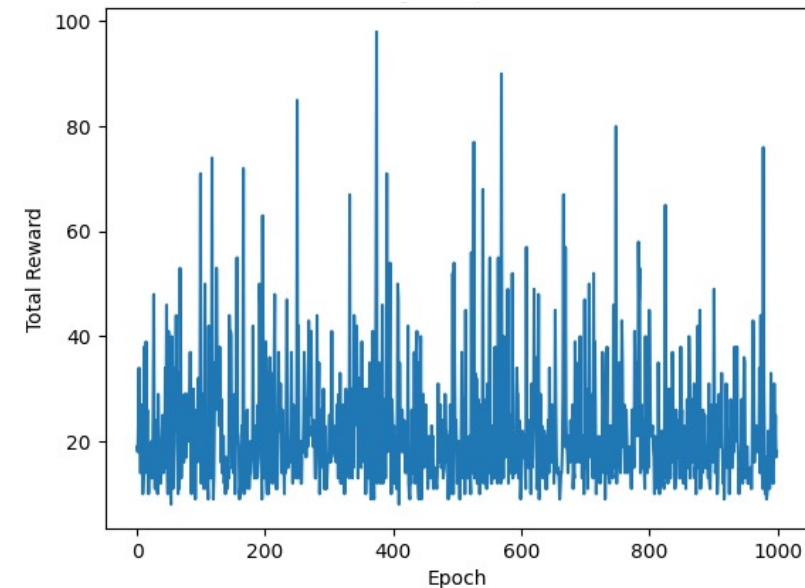
The model has done pretty well on the pendulum environment, which we can notice from the graph below that with quite a less of training, after 200 epochs the model has some significant progress and has stabilized near the end of the training phase with some positive rewards as well.



# Results

## PPO:

Proximal Policy Optimization algorithm has performance on our ForestFireEnv and CartPole have been suboptimal as they have not learnt the best policy with significant amount of training. This can be improved with better hyper parameter tuning, reward systems, Actor and Critic models.



# Key Observations

- Challenges we faced:

- Working in a MARL setup for different kinds of algorithms has been difficult to modify and exploring how the algorithms work has been fascinating.
- Modifying the environment and deep RL methods.
- Performance tuning to get the optimal results.

- Observations:

- Need a more comprehensive reward structure
- Most of the Deep RL methods have not been performing well on the stochastic setup as it is unpredictable.
- All the methods we have explored work a little better on existing environments with continuous action sets and less punishing reward systems.



# References

- Class notes and Piazza
- Assignment reference document
- Pytorch Documentation
- OpenAI Gym Documentation
- <https://github.com/msinto93/DDPG>
- <https://github.com/nikhilbarhate99/PPO-PyTorch>
- "Wildfire Detection using Deep Learning and Satellite Imagery" by J. He and H. Zhang. This paper proposes a deep learning approach to detecting wildfires using satellite imagery. The authors achieved an accuracy of 92% on their test set. DOI: 10.1109/JSTARS.2020.3043111
- Our CSE546 works on the Final project.



# CONTRIBUTION SUMMARY

Team Member	Project Part	Contribution (%)
tanwinch	ENV,DDPG,PPO,SAC,REPORT	33.33%
kumarraaj	ENV,DDPG,PPO,SAC,REPORT	33.33%
vkajeepe	DDPG,PPO,REPORT	33.33%

Though the work was assigned among teammates, each task has been done with the involvement of each teammate.

# THANK YOU!

