# CSE 4/546: REINFORCEMENT LEARNING

# Team 8: MULTI AGENT REINFORCEMENT LEARNING

Sri Amarnath Mutyala - : 50424956
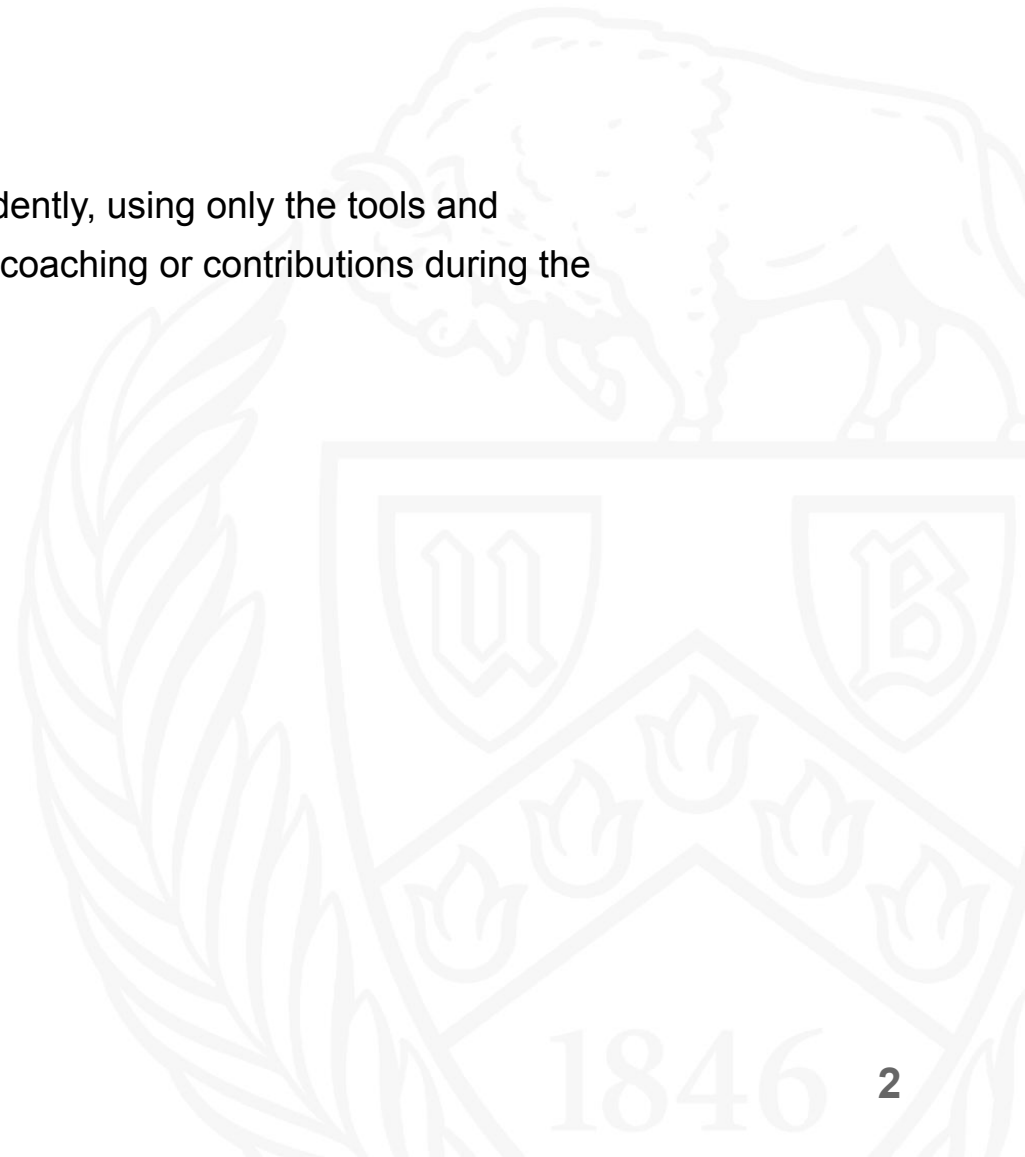
Naga Kiran Reddy Karnati - 50432242

Pritisman Kar - 50388920

University at Buffalo The State University of New York

# Academic Integrity Statement

"We certify that the code and data in this assignment were generated independently, using only the tools and resources defined in the course and that we did not receive any external help, coaching or contributions during the production of this work."

# Contents

- *Introduction*

    - Vehicle Scheduling Environment – Multi-Agent Grid World

    - Simple Adversary: OpenAI Multi-agent particle environment

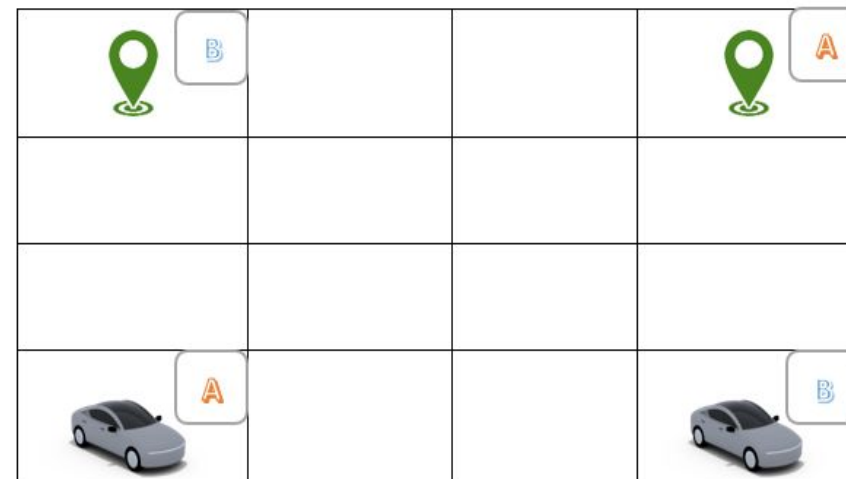- *Implementation, Results & Graphs*

    - Q learning on Vehicle Scheduling Environment – Multi-Agent Grid World (Independent Learning)

     - Multi Agent Deep Deterministic Policy Gradient (MADDPG) algorithm

    - MADDPG on Vehicle Scheduling Environment – Multi-Agent Grid World

    - Improved MADDPG on Vehicle Scheduling Environment – Multi-Agent Grid World

    - MADDPG on Vehicle Scheduling Environment – Simple Adversary

    - Comparison between Q learning and MADDPG algorithm on Vehicle Scheduling Environment - Multi-Agent Grid World

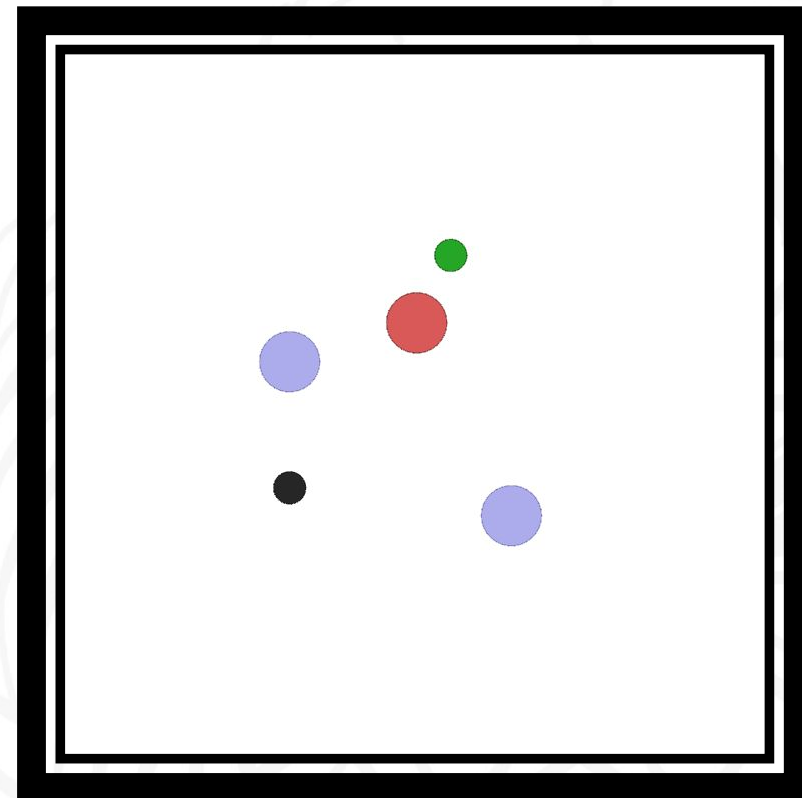- *Key observations*

- *References*

# Vehicle Scheduling Environment – Multi-Agent Grid World (MAGW)

- Two cars in a **4x4** environment

- 1st car – Goal - To reach top right of the environment

- 2nd car – Goal - To reach top left of the environment

- State space: **16 states: {s0, s1, s2,...s15}**

- Action space: {0: down, 1: up, 2: right, 3: left, 4: no move}

- Reward structure

    - Towards the target: **1**

    - Away from the target: **-3**
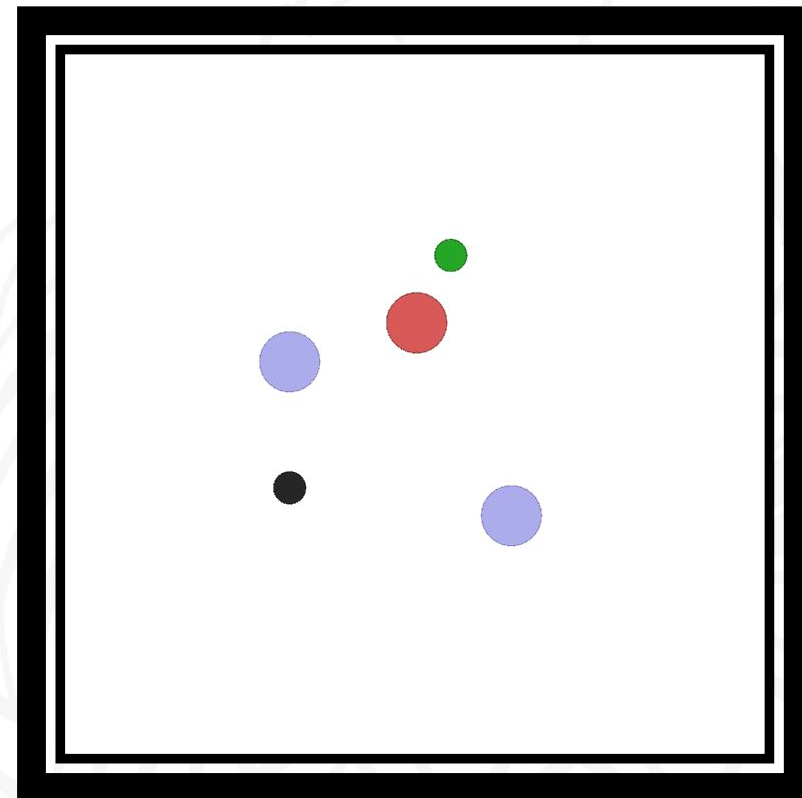
    - Stays in same position: **-5**

    - Reaches target: **100**



4

# Simple Adversary - OpenAI MA particle environment

- 3 agents – **1** adversary, **2** good agents (Physical deception)

- Environment – 2 landmark (**Green** – target landmark, **Black** – dummy landmark)

- Rewards:

  - For agents:

    - **Positive** reward - based on distance between the closest agent to target landmark

    - **Negative** reward – based on distance between the adversary to target landmark

  - For adversary:

    - **Positive reward** – based on distance between the adversary to target landmark



**https://www.pettingzoo.ml/mpe/mpe_simple_adversary.gif**

# Simple Adversary - Main features

| Parameter | Value |
|---|---|
| Actions | Discrete/Continuous |
| Agents | agents= [adversary_0, agent_0,agent_1] |
| Action Shape | (5) |
| Action Values | Discrete(5) |
| Observation Shape | (8) – adversary ,(10) – agents |
| Observation Values | (-inf,inf) |
| State Shape | (28,) [adversary + 2 agents] |
| State Values | (-inf,inf) |

# Q learning on Vehicle Scheduling Environment -MAGW

- Number of cars: **2**

- Episodes: **100000**

- Discount factor: **0.99**

- Learning rate: **0.15**

- Timesteps: **20**

# Results: Q-learning on MAGW

# Results: Q-learning on MAGW

# Multi Agent Deep Deterministic Policy Gradient

- First proposed by OpenAI in 2017

- Off-policy

- Algorithm which concurrently learns a Q-function and a policy

- Actor and Critic (i.e. and target networks for each) for each agent

- Can be used for environment with continuous action space and continuous environment

- Soft updates on parameters for each actor and critic network

- Uses Experience replay

- To ensure exploration, add Noise to deterministic policy gradient

- Issues – can overestimate the Q value in the critic network

**Multi-Agent Deep Deterministic Policy Gradient Algorithm**

For completeness, we provide the MADDPG algorithm below.

---

**Algorithm 1:** Multi-Agent Deep Deterministic Policy Gradient for $N$ agents

---

**for** episode $= 1$ to $M$ **do**

    Initialize a random process $\mathcal{N}$ for action exploration

    Receive initial state $\mathbf{x}$

    **for** $t = 1$ to max-episode-length **do**

        for each agent $i$, select action $a_i = \boldsymbol{\mu}_{\theta_i}(o_i) + \mathcal{N}_t$ w.r.t. the current policy and exploration

        Execute actions $a = (a_1, \ldots, a_N)$ and observe reward $r$ and new state $\mathbf{x}'$

        Store $(\mathbf{x}, a, r, \mathbf{x}')$ in replay buffer $\mathcal{D}$

        $\mathbf{x} \leftarrow \mathbf{x}'$

        **for** agent $i = 1$ to $N$ **do**

            Sample a random minibatch of $S$ samples $(\mathbf{x}^j, a^j, r^j, \mathbf{x}'^j)$ from $\mathcal{D}$

            Set $y^j = r_i^j + \gamma Q_i^{\boldsymbol{\mu}'}(\mathbf{x}'^j, a_1', \ldots, a_N')|_{a_k' = \boldsymbol{\mu}_k'(o_k^j)}$

            Update critic by minimizing the loss $\mathcal{L}(\theta_i) = \frac{1}{S}\sum_j \left(y^j - Q_i^{\boldsymbol{\mu}}(\mathbf{x}^j, a_1^j, \ldots, a_N^j)\right)^2$

            Update actor using the sampled policy gradient:

$$\nabla_{\theta_i} J \approx \frac{1}{S}\sum_j \nabla_{\theta_i} \boldsymbol{\mu}_i(o_i^j) \nabla_{a_i} Q_i^{\boldsymbol{\mu}}(\mathbf{x}^j, a_1^j, \ldots, a_i, \ldots, a_N^j)|_{a_i = \boldsymbol{\mu}_i(o_i^j)}$$

        **end for**

        Update target network parameters for each agent $i$:

$$\theta_i' \leftarrow \tau\theta_i + (1 - \tau)\theta_i'$$

    **end for**

**end for**

---

# Implementation details:

**Every Agent has**

Actor Network:

    Inputs: **States, actions**
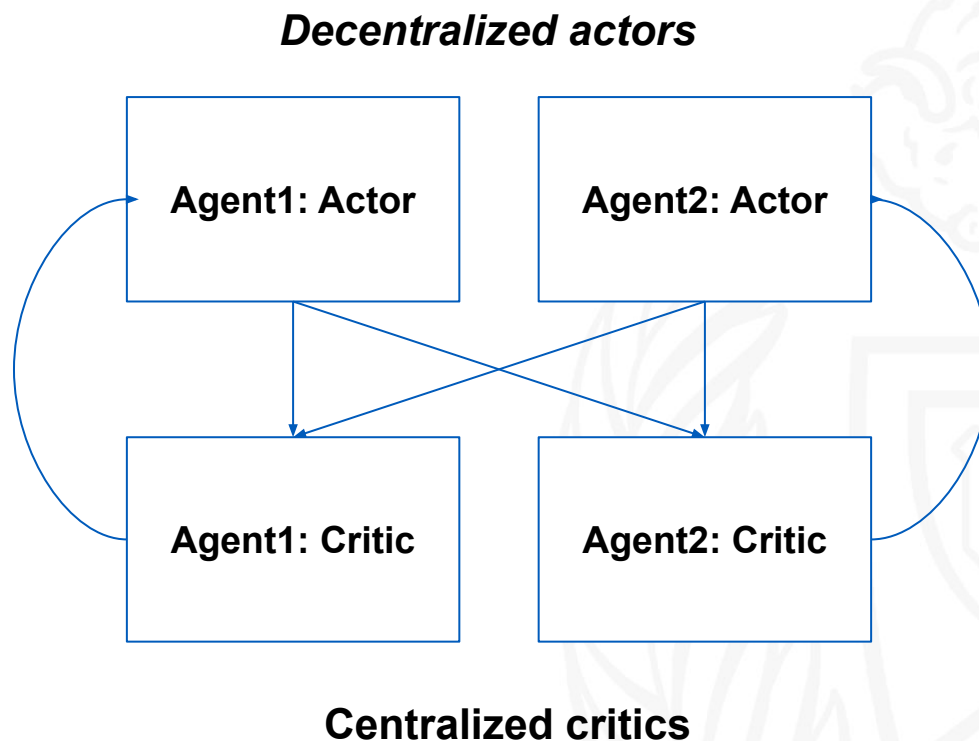
    Outputs: **Probs**

Critic Network:

    Inputs: **states, actions**

    Outputs: **Q values**

**To freeze weights to avoid running targets**

Target Actor Network (i.e. performed soft updates)
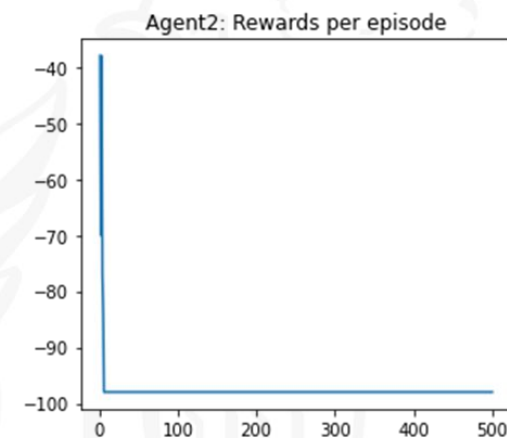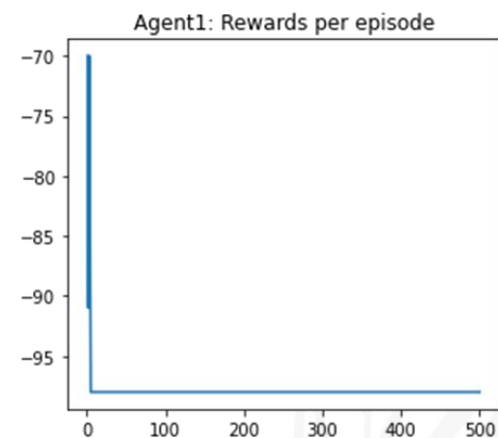
Target Critic Network (i.e. performed soft updates)

*Decentralized actors*

| Agent1: Actor | Agent2: Actor |
|---|---|
| Agent1: Critic | Agent2: Critic |

**Centralized critics**

12

# MADDPG on MAGW

- Summary of actor network –

  - **2 hidden layers (64 each), softmax output, learning rate - alpha: 0.01**

- Summary of the critic network –

  - **2 hidden layers(64 each), relu output, learning rate - beta: 0.01**

- Number of episodes – **500**

- Gamma - **0.95**

- Batch size - **64**

# Results: MADDPG on MAGW

# Improved MADDPG on MAGW

- **ε-greedy approach** even after applying noise to actions chosen from deterministic policy

- Number of episodes – **1000**

- Batch size: **128**

- Actor network: alpha: **0.001**

- Critic network: beta: **0.001**

# Results: Improved MADDPG on MAGW

# MADDPG on Simple Adversary - Implementation Details

- Summary of actor network –

    - **2 hidden layers (64 each), softmax output, learning rate - alpha: 0.01**

- Summary of the critic network –

    - **2 hidden layers(64 each), relu output, learning rate - beta: 0.01**

- Number of episodes – **15000**

- Gamma - **0.95**

- Batch size - **1024**

- Maximum timesteps: **25**

- Learning: **Every 100 steps**

# Results – MADDPG on Simple Adversary

# Comparison of Q-learning & Improved MADDPG on MAGW

**Q-Learning**

**Improved MADDPG**

# Key observations

- As noticed in the previous slide, the Q learning is not working well for the Vehicle Scheduling environment.

- The MADDPG algorithm is working better when compared to Q learning algorithm.

- Proper attention should be given while implementing the MADDPG algorithm since it may lead to over-estimation of the Q-value using the Critic network.

- MADDPG is working well for continuous action-state value environment (Simple Adversary)

# Other Trails & Extensions

- We were successful in setting up the Mujoco environment in CCR.

- We have also run benchmarks on Ant-v2 and HalfCheetah-v2 environments using RLLib package from the Ray Team.

- We eventually wanted to write our algorithm in the Ray tune framework as they have software primitives in RL and distributed computing.

- We have also defined a fleet scheduling environment with good reasonable assumptions.

- TODO: Apply the MADDPG algorithm to the fleet scheduling environment.

# Contribution

| Name | Contribution |
|---|---|
| Naga Kiran Reddy Karnati | 33% |
| Sri Amarnath Mutyala | 33% |
| Pritisman Kar | 33% |

# Trello workspace timeline

# Trello workspace timeline

# Trello workspace timeline

# Reference

- [1706.02275.pdf (arxiv.org)](https://arxiv.org)

- [GitHub - philtabor/Multi-Agent-Deep-Deterministic-Policy-Gradients: A Pytorch implementation of the multi agent deep deterministic policy gradients (MADDPG) algorithm](https://github.com)

- [https://www.pettingzoo.ml/mpe/simple_adversary.html](https://www.pettingzoo.ml/mpe/simple_adversary.html)

- [Efficient Large-Scale Fleet Management via Multi-Agent Deep Reinforcement Learning | Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining](https://dl.acm.org)

THANK YOU