# A Reinforcement Learning Approach to Ship Towing Using Two Tugboats

Name: Mohammad Saifullah Khan
Roll No.: 21169

TEAM:
Rahul Kulkarni
Mohammad Saifullah Khan
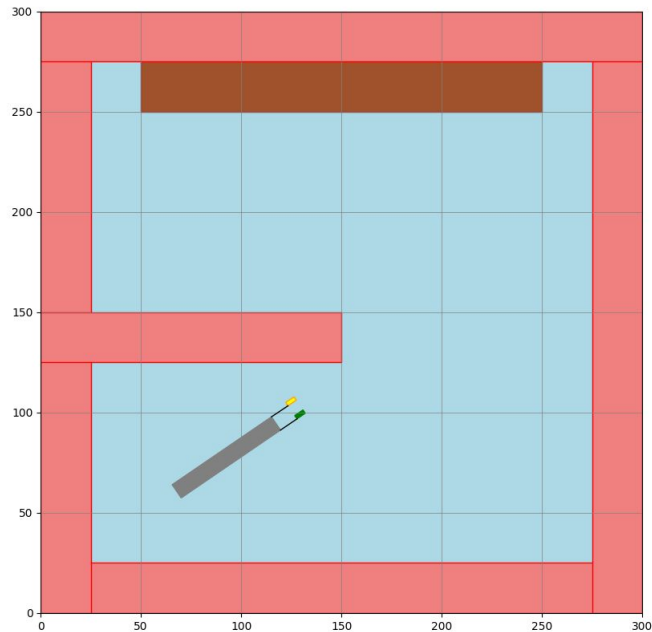
# Problem Statement

Towing a ship to target location by two tugboats using Multi-Agent Deep Deterministic Policy Gradient (MADDPG)

# Challenges

- Coordination among agents (tugboats).
- High dimensional and continuous action space.
- Scalability and complexity of task
- Reward shaping and credit assignment.
- Adaptation to different ship types and conditions.
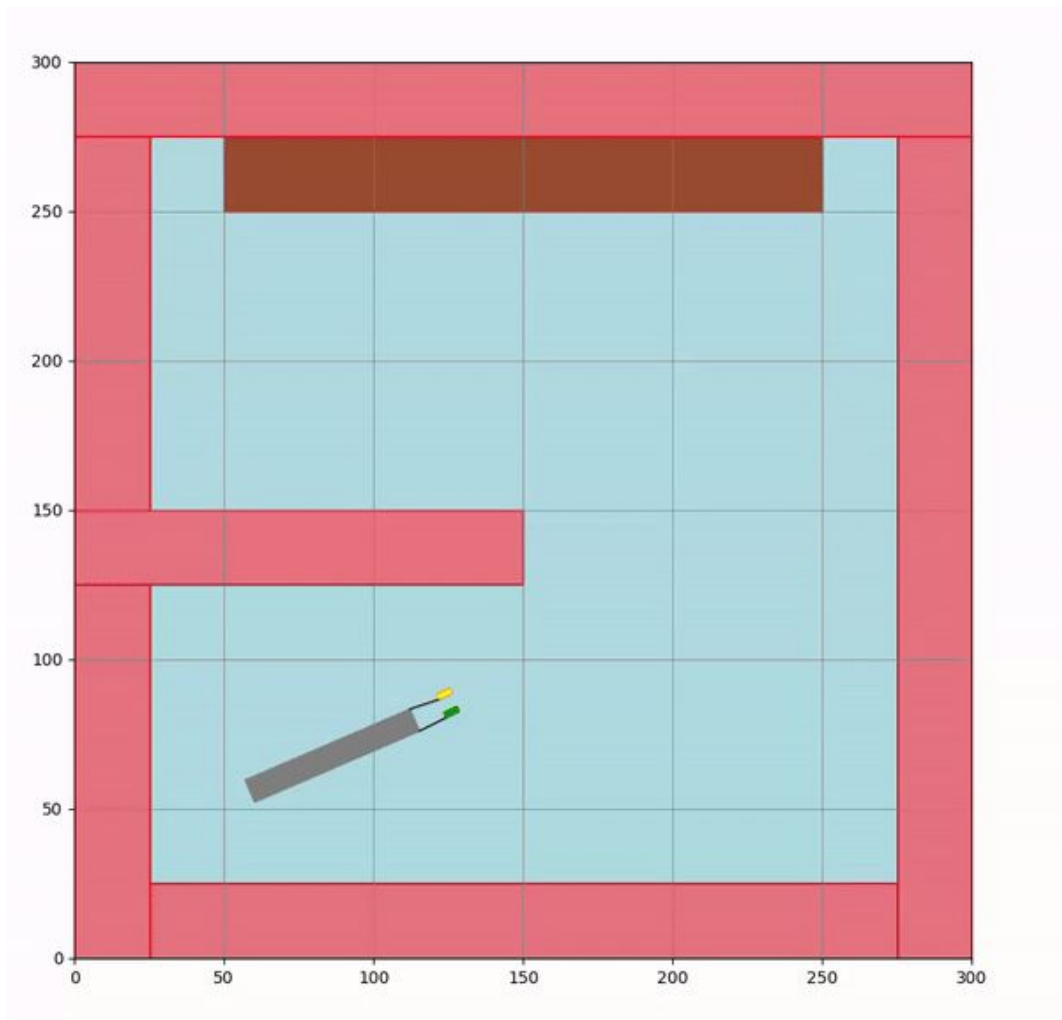- Real world transfer of simulation.

# Approach



Agents = tugboats.
Obstacles: red coloured.
Dock: brown coloured.
Water: blue coloured

Built custom environment.
- Observation space of each tugboat: ship position and orientation, position and orientation of own and other agent, distance of ship from target (Euclidean), rope length.
- Action space of each tugboat: velocity in x and y direction.
- Rope (black coloured) considered as a stiff rod for simplicity.
- A force applied on ship due to movement of tugboats.

## Reward structure:

$$R_{to\_target} = - ds/(grid\_size*0.004)$$

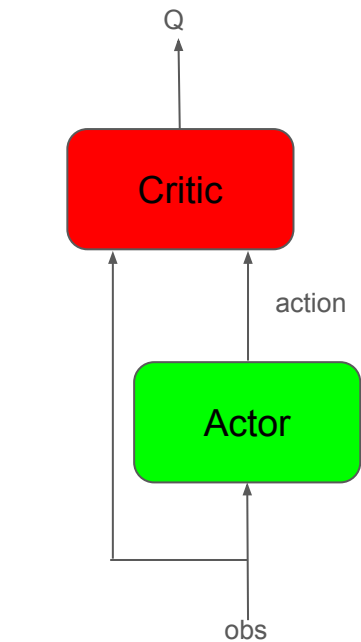$$R_{proximity} = \text{Penalty for coming close to obstacles}$$
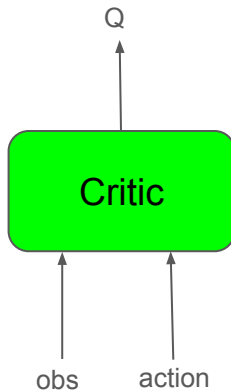$$\text{(for both ship and tugboats)}$$

$$R_{in\_target} = 1000$$

$$R_{rope\_length} = \text{penalty for exceeding rope length}$$

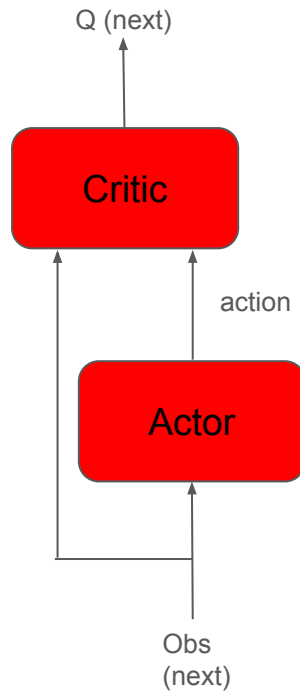$$R_{total} = R_{to\_target} + R_{proximity} + R_{rope\_length} + R_{in\_target}$$

# DDPG



$$\text{minimize}(Q - |\ \text{reward} + \text{discount} * Q_{next}|)$$

# MADDPG

- Soft update of parameters.
  $$\theta_{target} \leftarrow \tau\theta_{online} + (1-\tau)\theta_{target}$$

- Handles continuous action spaces.
- Sample efficiency
- Training stability.
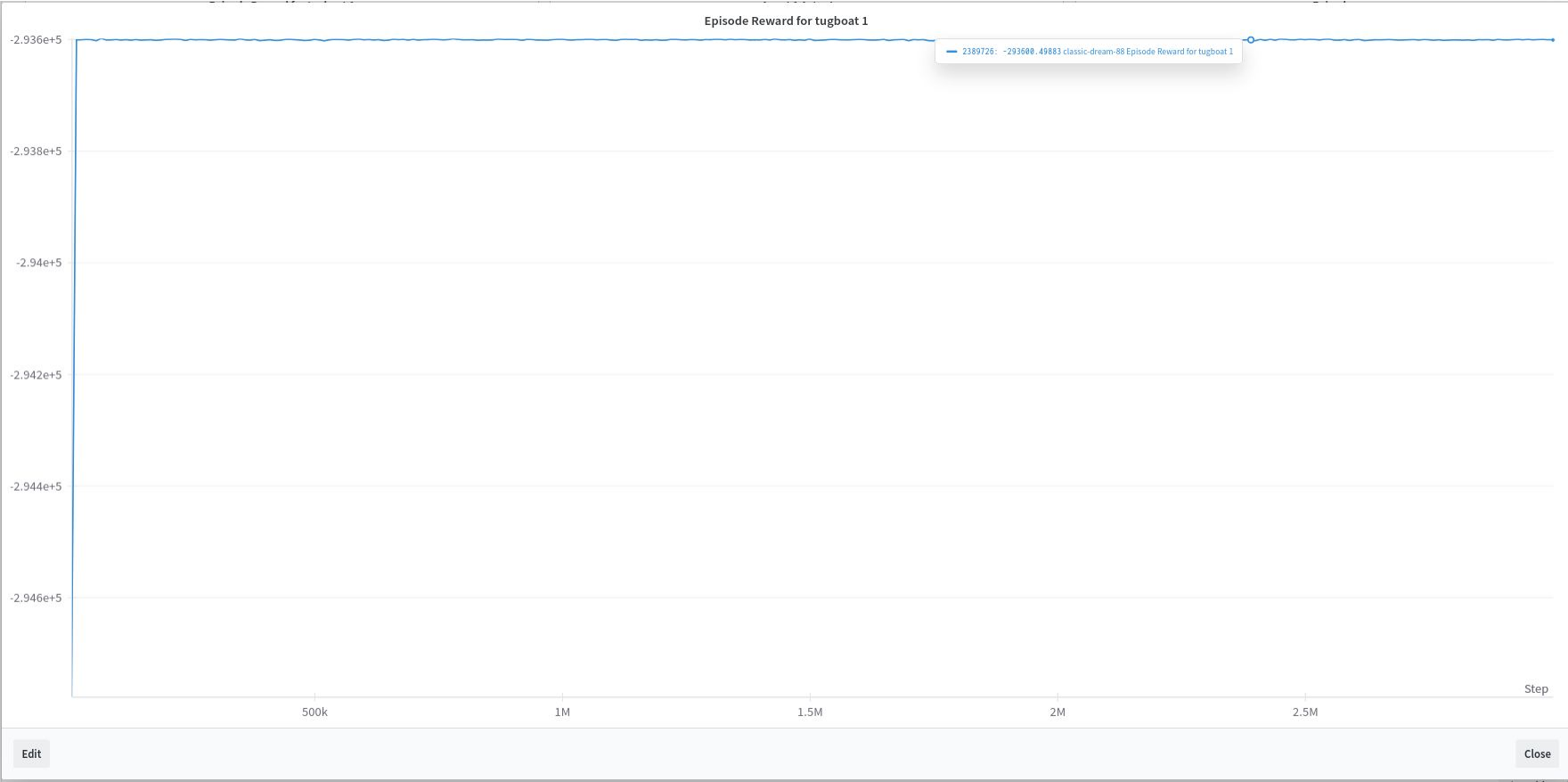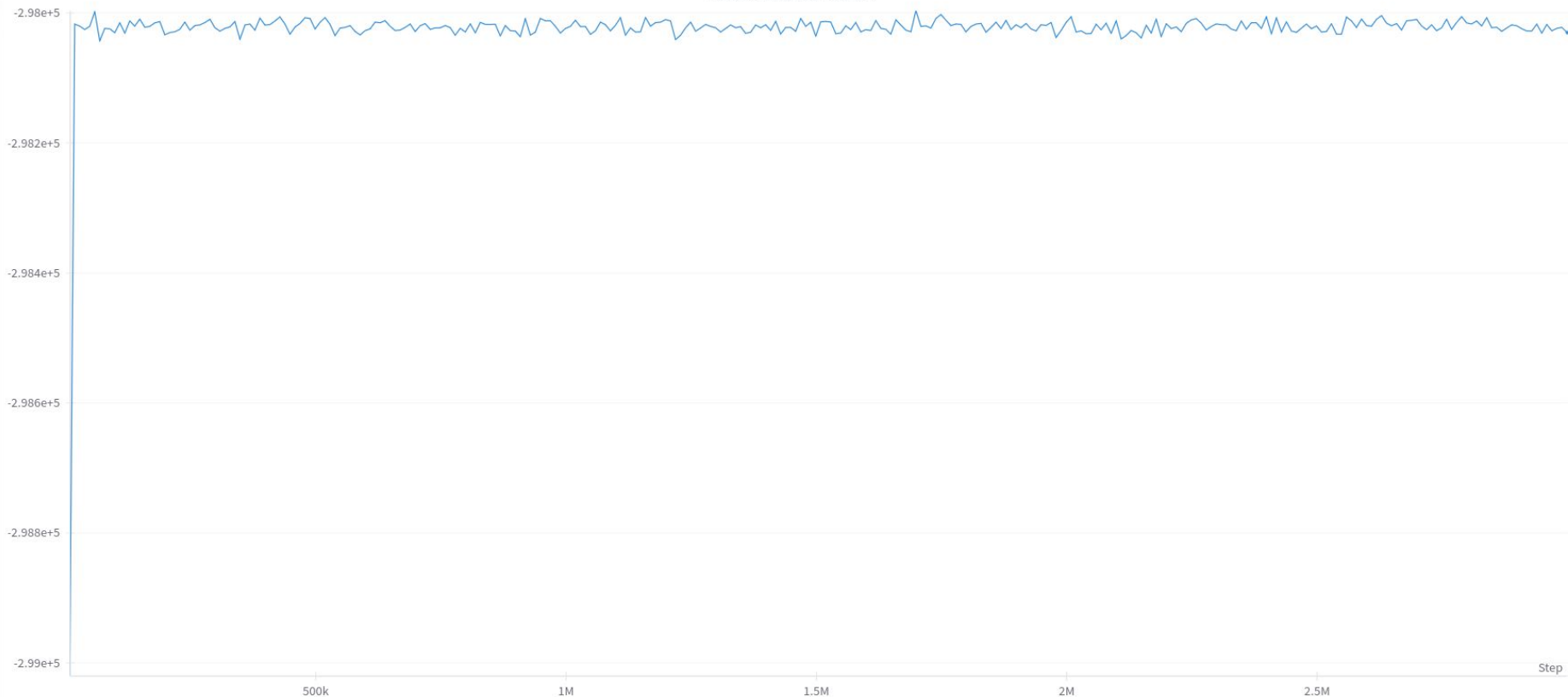- Centralized training, decentralized execution.



Source: https://arxiv.org/pdf/1706.02275

# Parameters

- Gamma = 0.99
- Adam optimizer, learning rate = 0.003.
- Tau = 0.01
- Training episodes = 300, each with 2500 steps.

# Results



Episode Reward for tugboat 1

2389726:  -293600.49883 classic-dream-88 Episode Reward for tugboat 1

Edit                                                                                                                    Close
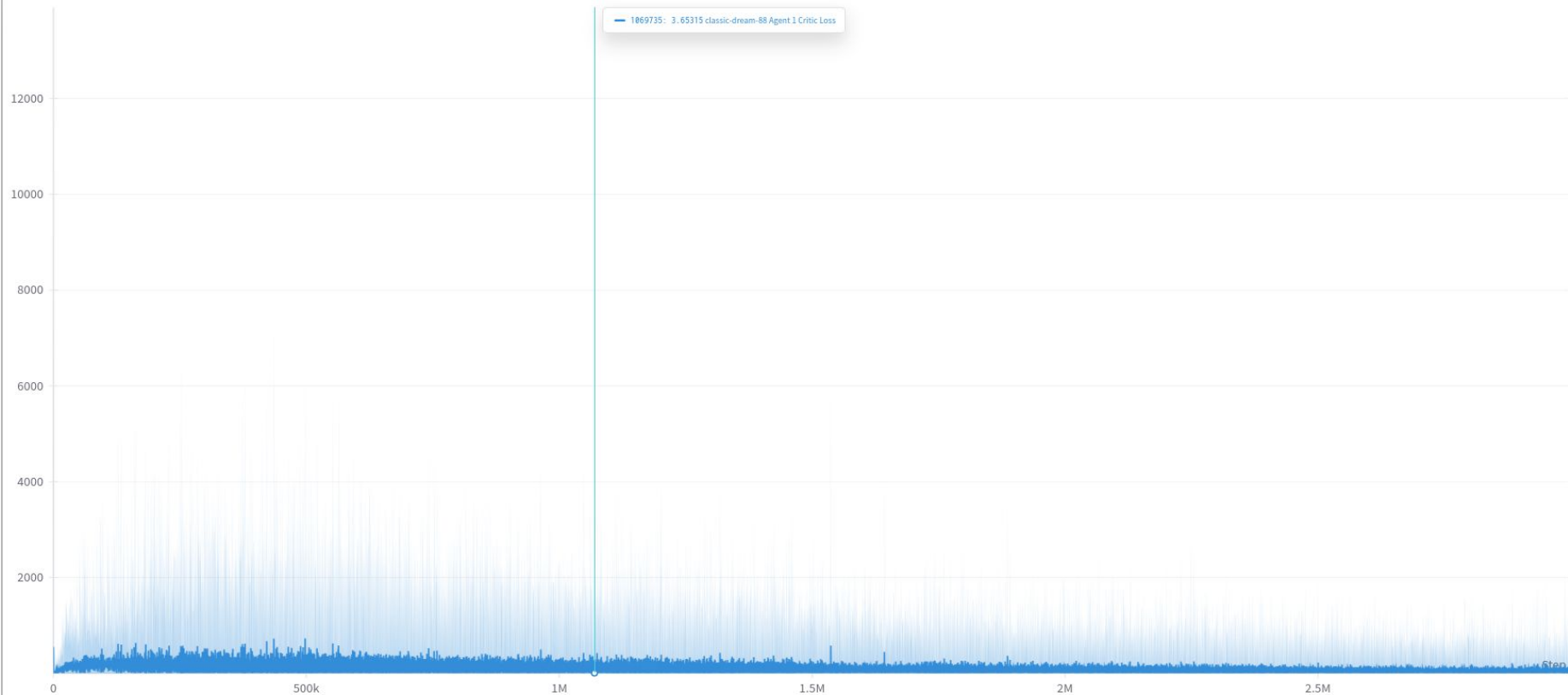
Episode Reward for tugboat 2

Edit                                                                                    Close

# Agent 1 Critic Loss



— 1069735: 3.65315 classic-dream-88 Agent 1 Critic Loss

Edit

Close

# Agent 2 Critic Loss

# Agent 1 Actor Loss

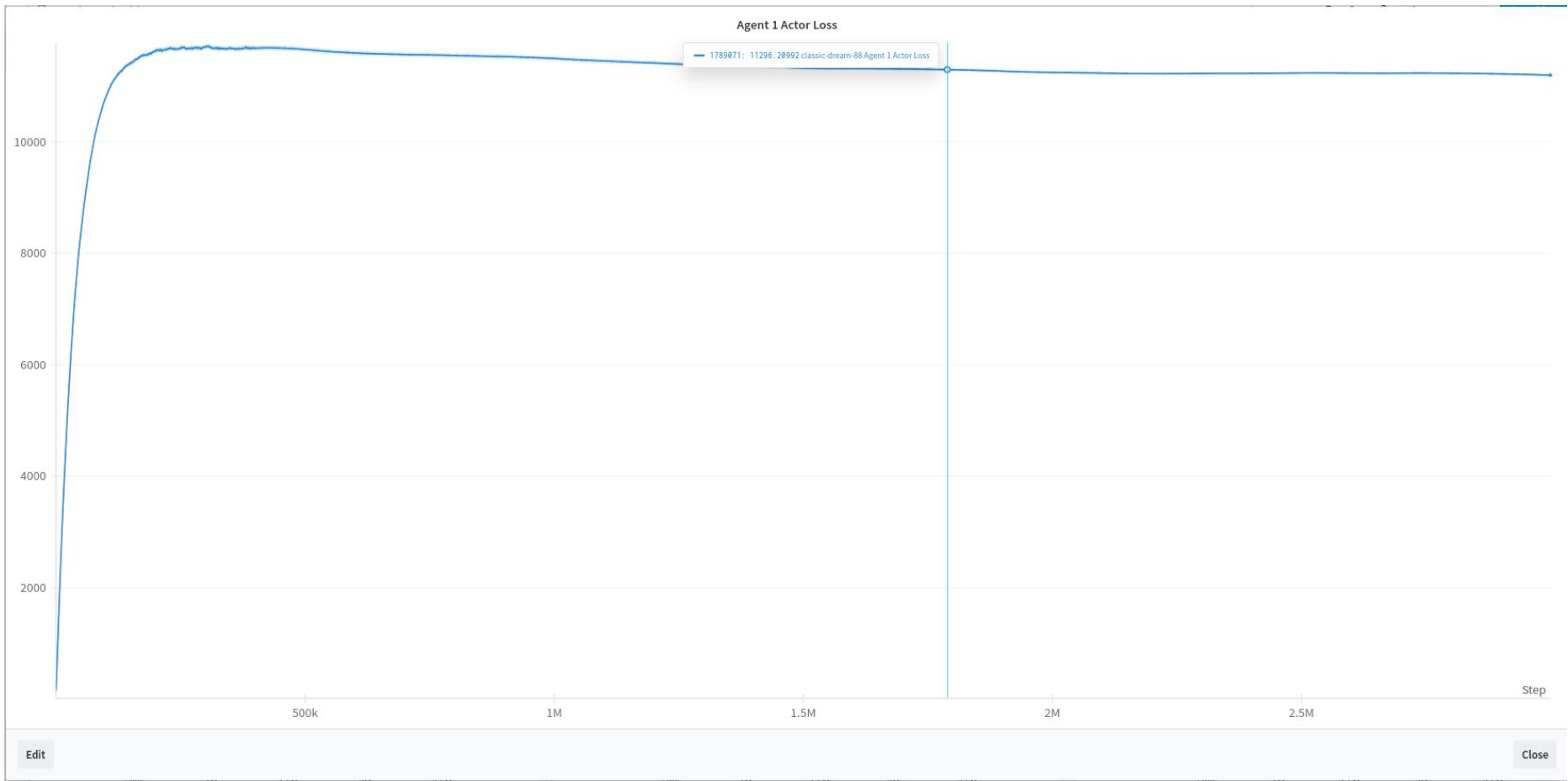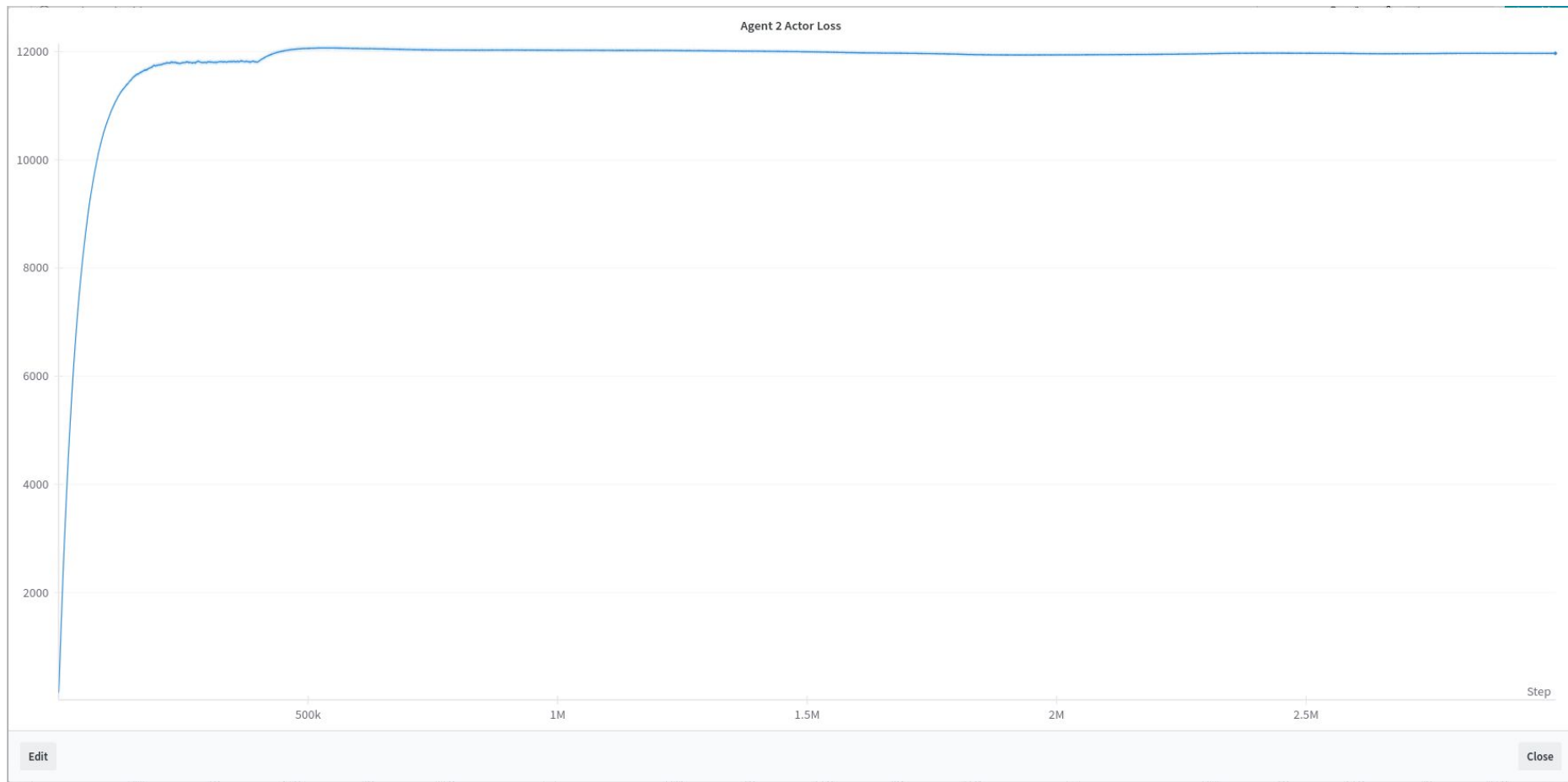— 1789071: 11296.20992 classic-dream-88 Agent 1 Actor Loss

Agent 2 Actor Loss

# For random steps



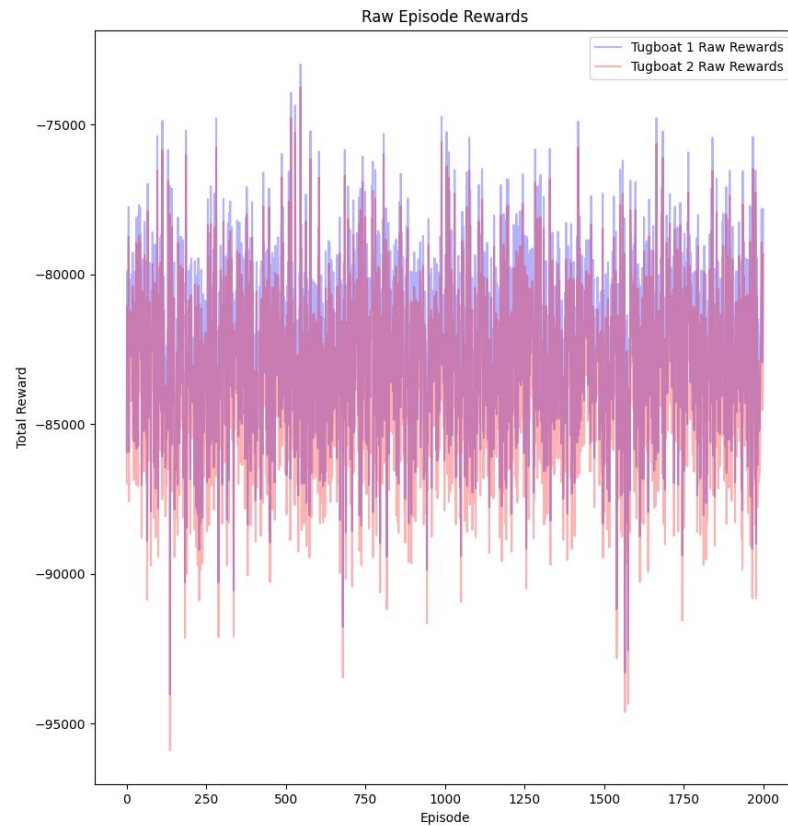Raw Episode Rewards

# Future Work

- Improve reward structure.
- Random spawn in the grid.
- Add dynamic obstacles.
- Managing traffic of multiple tugboat – ship combinations.