

SwarmAI

20190196 TaeHyung Kim

20190321 Seunghwan Song

20190456 SangHyun Lee

Intro

- Complex interactions
- Behaviors for survival
- Model ecosystems



Intro

- Swarm behavior
- Disturbing predator
- Increase chances of survival



Background

- Fish agent (Discrete Action Multi Agent Actor-Critic)

Algorithm 1 Discrete Action Multi-Agent Actor Critic for N agents

```
1: for episode = 1 to  $M$  do
2:   Receive initial state  $x$ 
3:   for  $t = 1$  to max-episode-length do
4:     for each agent  $i$  do
5:       Select action  $a_i \sim \mu_{\theta_i}(o_i)$  w.r.t. the current policy and exploration
6:     end for
7:     Execute actions  $\mathbf{a} = (a_1, \dots, a_N)$  and observe reward  $r$  and new state  $x'$ 
8:     Store  $(x, a, r, x')$  in replay buffer  $\mathcal{D}$ 
9:      $x \leftarrow x'$ 
10:  end for
11:  for agent  $i = 1$  to  $N$  do
12:    Sample a random minibatch of  $S$  samples  $(x^j, a^j, r^j, x'^j)$  from  $\mathcal{D}$ 
13:    Set  $y^j = r^j + \gamma Q_{\theta'_i}^i(x', a_1, \dots, a_N)|_{a_k = \mu_{\theta_k}(o_k)}$ 
14:    Update critic by minimizing the loss
```

$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_j \left(y^j - Q_{\theta_i}^i(x^j, a_1^j, \dots, a_N^j) \right)^2$$

```
15:    Update actor using the sampled policy gradient
```

$$\nabla_{\theta_i} J \approx \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_{\theta_i}(o_i^j) \nabla_{a_i} Q_{\theta_i}^i(x^j, a_1^j, \dots, a_N^j)|_{a_i = \mu_{\theta_i}(o_i^j)}$$

```
16:  end for
17:  Update target network parameters for each agent  $i$ :
```

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

```
18: end for
```

Background

- Fish agent (Discrete Action Multi Agent Actor-Critic)

Algorithm 1 Discrete Action Multi-Agent Actor Critic for N agents

```
1: for episode = 1 to  $M$  do
2:   Receive initial state  $x$ 
3:   for  $t = 1$  to max-episode-length do
4:     for each agent  $i$  do
5:       Select action  $a_i \sim \mu_{\theta_i}(o_i)$  w.r.t. the current policy and exploration
6:     end for
7:     Execute actions  $\mathbf{a} = (a_1, \dots, a_N)$  and observe reward  $r$  and new state  $x'$ 
8:     Store  $(x, \mathbf{a}, r, x')$  in replay buffer  $\mathcal{D}$ 
9:      $x \leftarrow x'$ 
10:  end for
11:  for agent  $i = 1$  to  $N$  do
12:    Sample a random minibatch of  $S$  samples  $(x^j, a^j, r^j, x'^j)$  from  $\mathcal{D}$ 
13:    Set  $y^j = r^j + \gamma Q_{\theta_i'}^i(x', a_1, \dots, a_N)|_{a_k = \mu_{\theta_k}(o_k)}$ 
14:    Update critic by minimizing the loss
```

$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_j \left(y^j - Q_{\theta_i}^i(x^j, a_1^j, \dots, a_N^j) \right)^2$$

```
15:   Update actor using the sampled policy gradient
```

$$\nabla_{\theta_i} J \approx \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_{\theta_i}(o_i^j) \nabla_{a_i} Q_{\theta_i}^i(x^j, a_1^j, \dots, a_N^j)|_{a_i = \mu_{\theta_i}(o_i^j)}$$

```
16:   end for
17:   Update target network parameters for each agent  $i$ :
```

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

```
18: end for
```

Discrete Action select



Background

- Fish agent (Discrete Action Multi Agent Actor-Critic)

Algorithm 1 Discrete Action Multi-Agent Actor Critic for N agents

```
1: for episode = 1 to  $M$  do
2:   Receive initial state  $x$ 
3:   for  $t = 1$  to max-episode-length do
4:     for each agent  $i$  do
5:       Select action  $a_i \sim \mu_{\theta_i}(o_i)$  w.r.t. the current policy and exploration
6:     end for
7:     Execute actions  $\mathbf{a} = (a_1, \dots, a_N)$  and observe reward  $r$  and new state  $x'$ 
8:     Store  $(x, a, r, x')$  in replay buffer  $\mathcal{D}$ 
9:      $x \leftarrow x'$ 
10:  end for
11:  for agent  $i = 1$  to  $N$  do
12:    Sample a random minibatch of  $S$  samples  $(x^j, a^j, r^j, x'^j)$  from  $\mathcal{D}$ 
13:    Set  $y^j = r^j + \gamma Q_{\theta_i'}^i(x', a_1, \dots, a_N)|_{a_k = \mu_{\theta_k}(o_k)}$ 
14:    Update critic by minimizing the loss
15:    Update actor using the sampled policy gradient
16:  end for
17:  Update target network parameters for each agent  $i$ :
18:   $\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$ 
```

Select action $a_i \sim \mu_{\theta_i}(o_i)$ w.r.t. the current policy and exploration

Discrete Action select

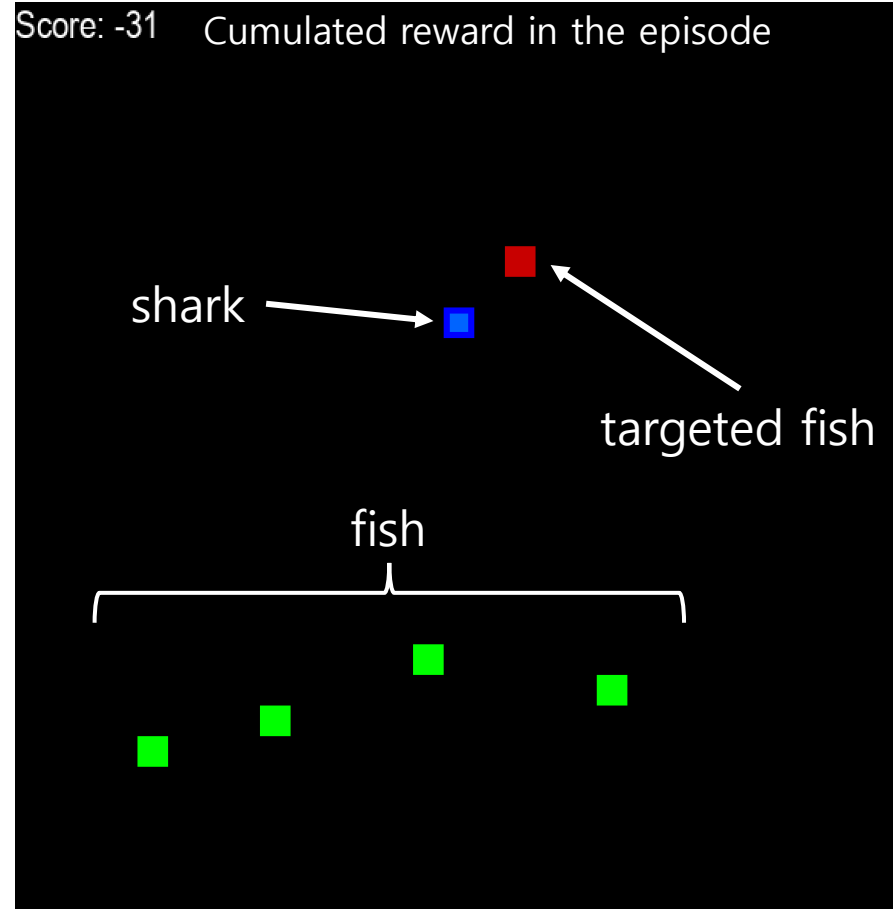
$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_j (y^j - Q_{\theta_i}^i(x^j, a_1^j, \dots, a_N^j))^2$$

Using Softmax

$$\nabla_{\theta_i} J \approx \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_{\theta_i}(o_i^j) \nabla_{a_i} Q_{\theta_i}^i(x^j, a_1^j, \dots, a_N^j)|_{a_i = \mu_{\theta_i}(o_i^j)}$$

Method

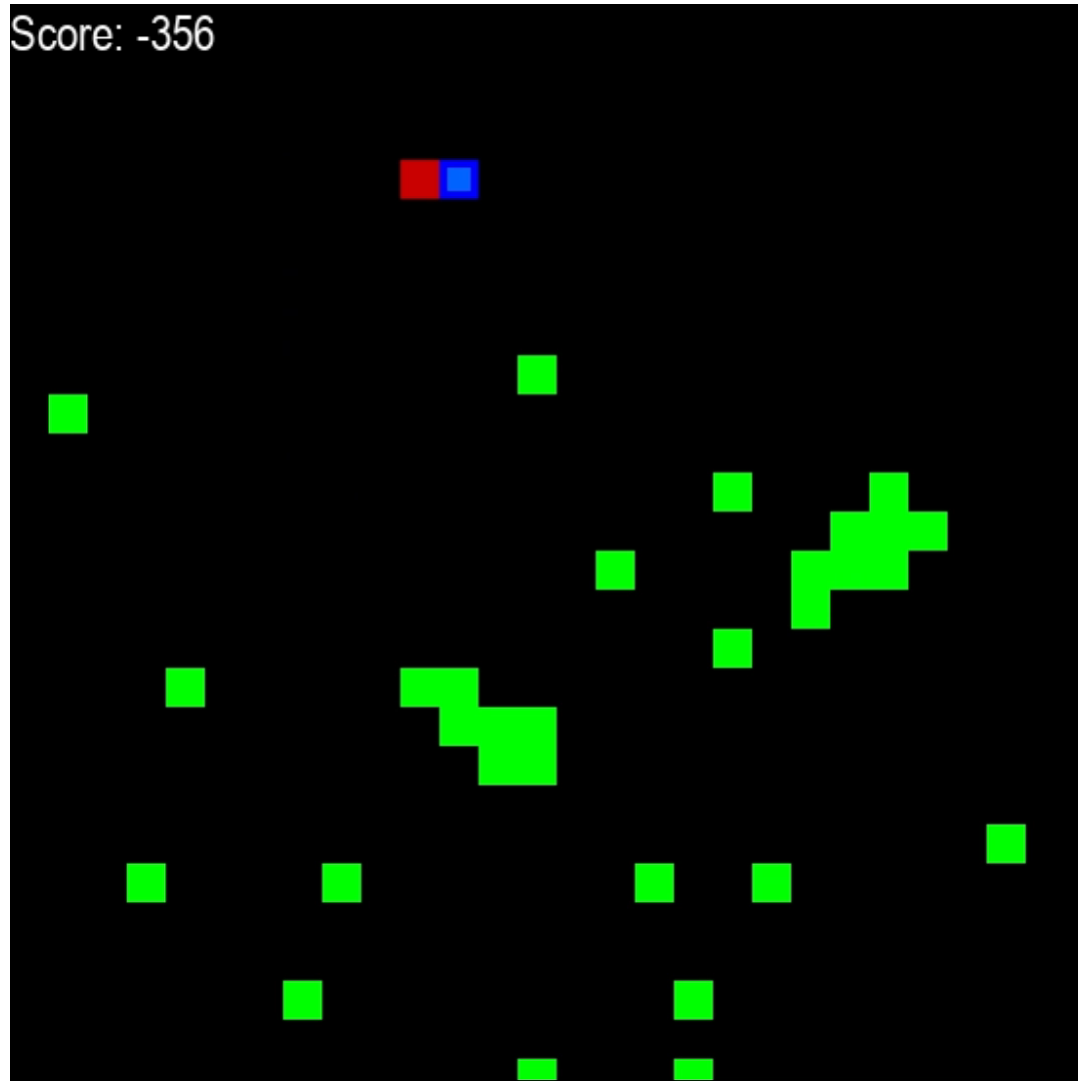
- Simulation in 2D pixel space
- Periodic boundary
- Partially observable



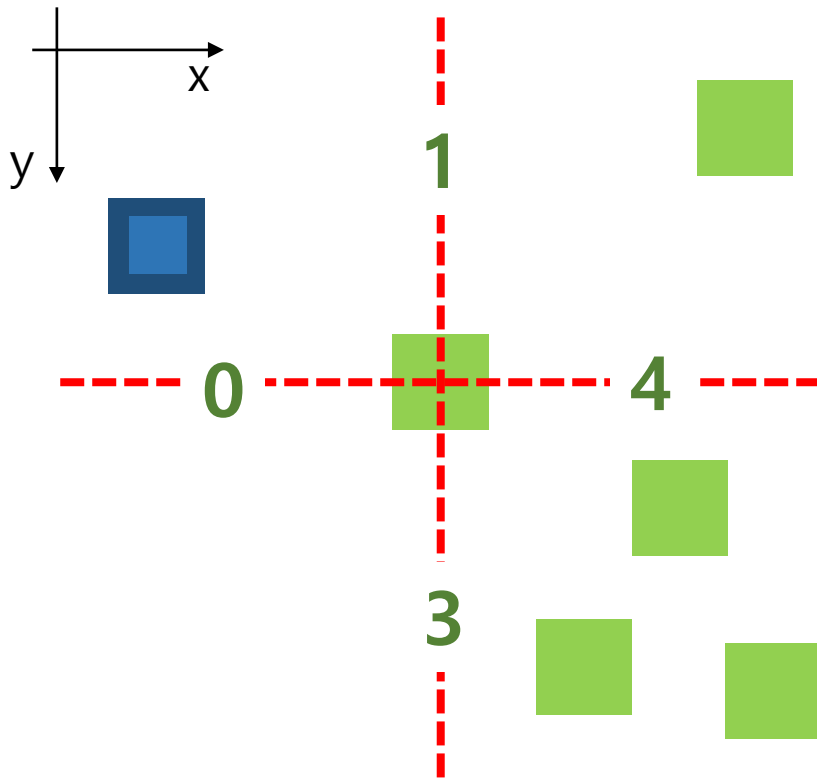
Method

- Game loop (2000 steps per episode)
 1. Fish moves w.r.t policy
 2. Shark moves
 3. Check fish eaten => -1 **reward** per fish eaten
 4. Reset shark target to closest fish

Example run

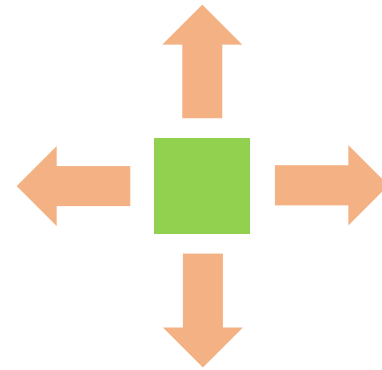


State space



- Number of fish in each direction: (1,3,0,4)
- Shark direction: (1,0,1,0) w.r.t. (up, down, left, right)

Action space



Go one block among four directions: up down left right

Method

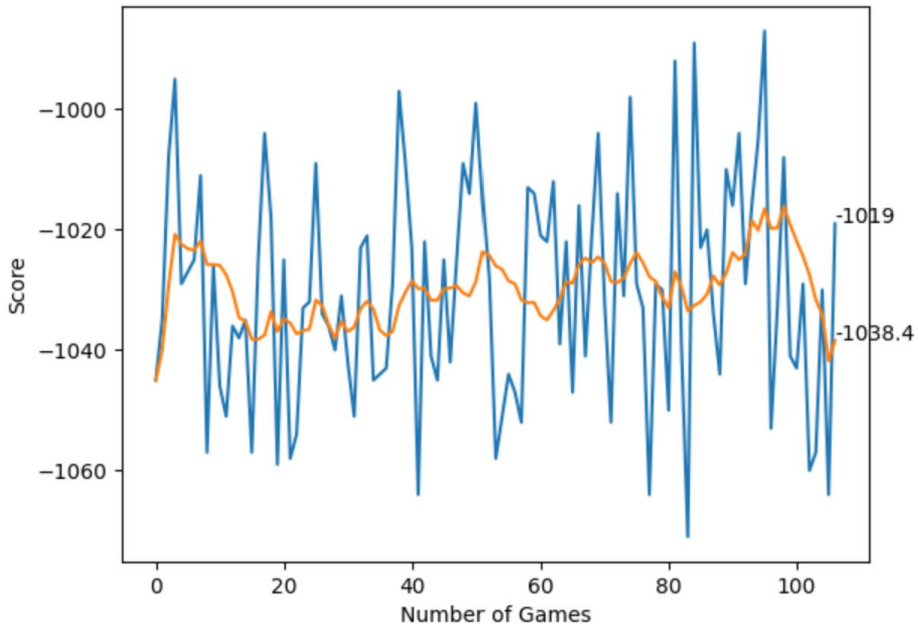
- Shark (challenge)

Algorithm Measuring the swarm size of a fish by the predator

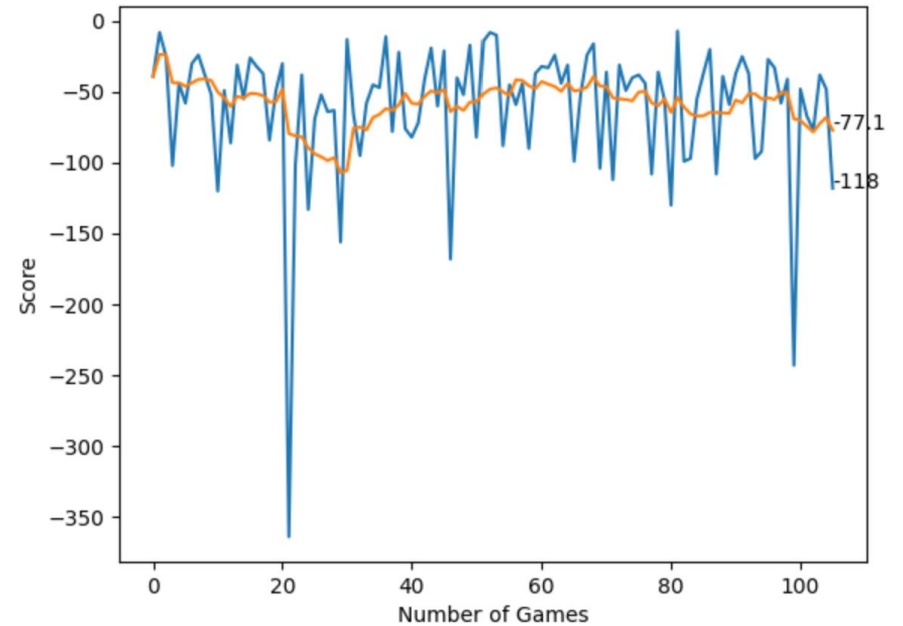
```
1: Get the target fish and its x,y coordinates  $x, y$ 
2: Initialize the  $swarmsize = 0$ 
3: for each  $fish$  do
4:   if  $fish$  is dead or  $fish$  is current target fish then
5:     continue
6:   else
7:     Take care of periodic boundary as follows
8:      $dx = \min(abs(x - fish.x), WIDTH - abs(x - fish.x))$ 
9:      $dy = \min(abs(y - fish.y), HEIGHT - abs(y - fish.y))$ 
10:     $distance = \sqrt{dx * dx + dy * dy}$ 
11:    if distance  $\leq$  SWARM RADIUS then
12:       $swarmsize += 1$ 
13:    end if
14:  end if
15: end for
16: return  $swarmsize + 1$ 
```

Baselines

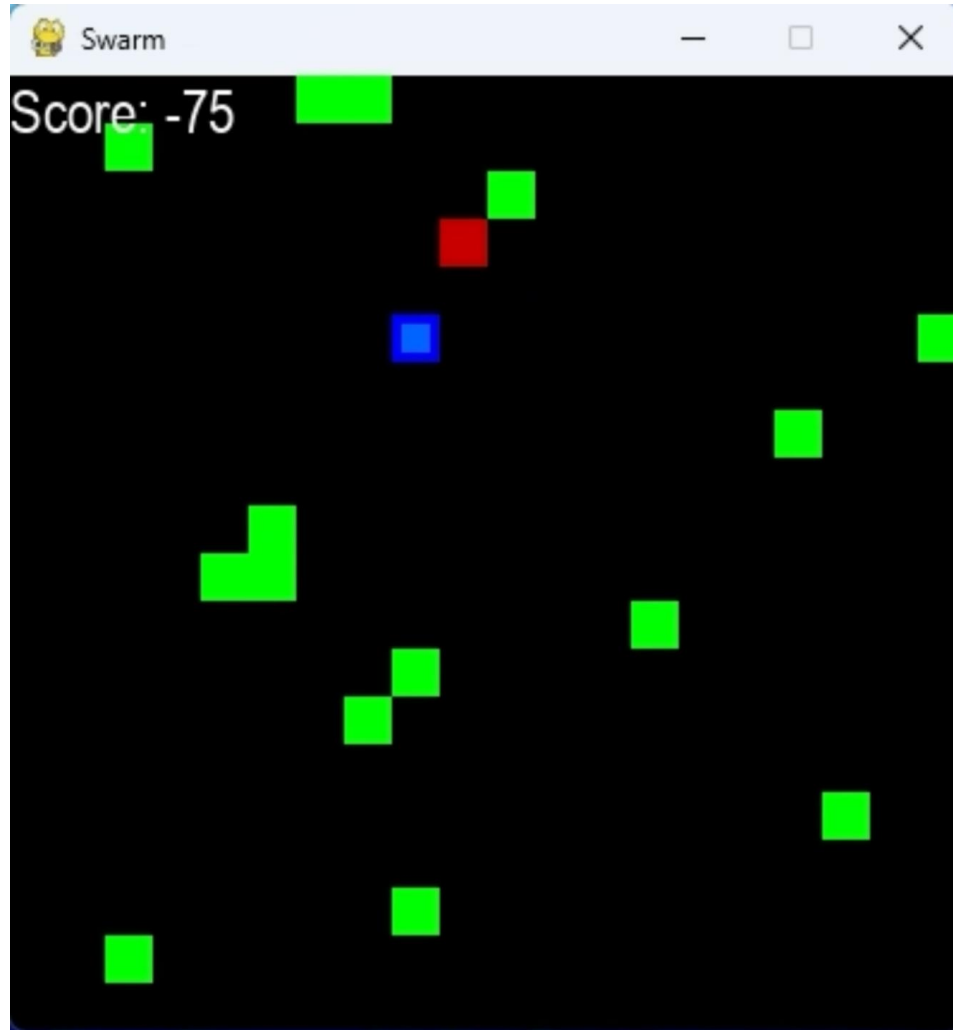
random



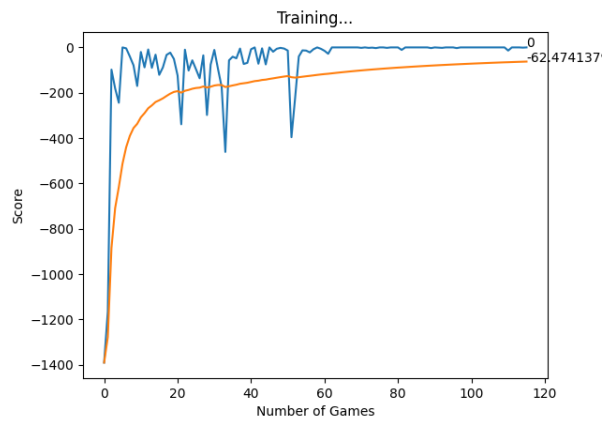
One direction



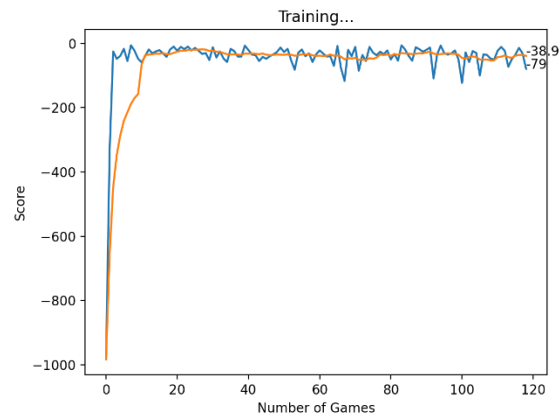
Experiment



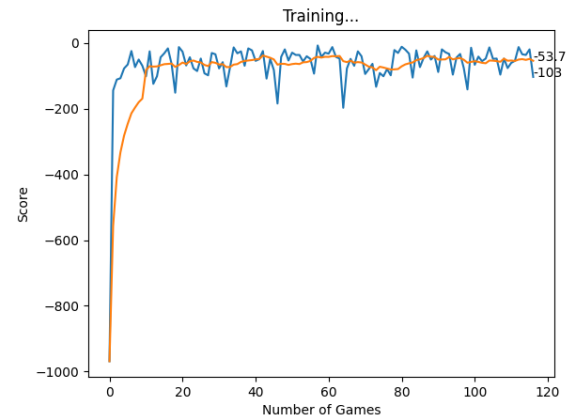
Result



200 x 200



400 x 400



600 x 600

Discussion

- **Environmental Variations and Robustness**
 - Effective group behavior across different environmental conditions.
- **Behavioral Observations**
 - Tend to form cohesive groups
 - Some fish remained alone and moved in a single direction



Conclusion

- Potential of MADDPG in **predator-prey** ecosystem
- Future works will explore more complex and realistic environments.
 - Refining the model to minimize human intervention.



Thank you