

Learning Complex Group Behaviours in a Multi-Agent Competitive Environment

Kazım Sanlav

Department of Industrial Engineering

Bilkent University

Afshan Nabi

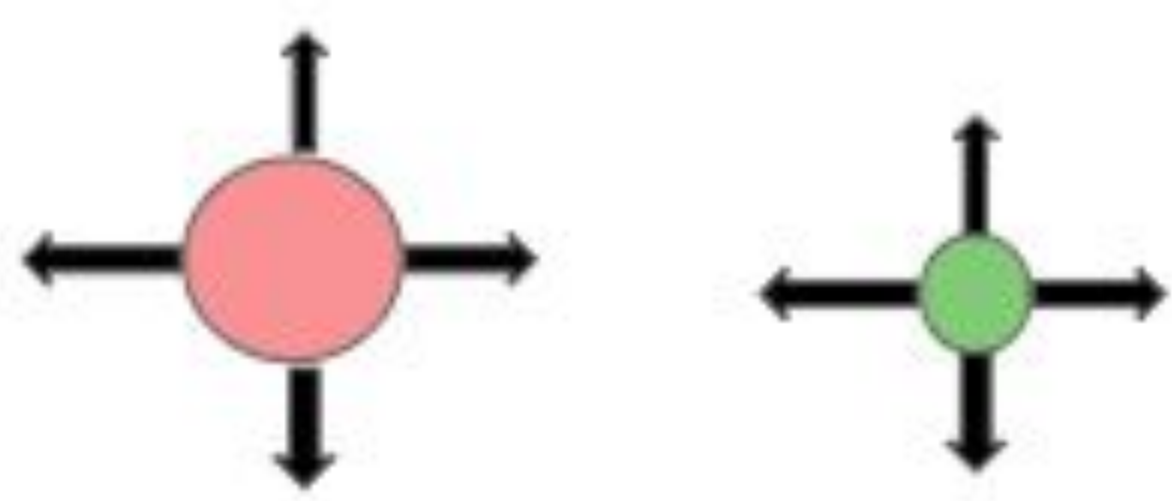
Department of Molecular Biology and Genetics

Introduction

The aim of this project is to use reinforcement learning algorithms to simulate predator-prey dynamics, such as those seen when predator fish and prey fish interact in nature. In order to do so, we have adapted OpenAI's multiagent-particle-envs to simulate two competing groups of fish- a predator group (has a single predator) and a prey group (2 or 10 prey). We used 2 different Q-learning algorithms (Q-Network with Experience Replay & Q-Network with Experience Replay and Fixed Q-Target) and a Monte-Carlo Policy Gradient algorithm to train these predator and prey agents. The predator is rewarded based on the number of times they collide with ("eat") prey. The prey are rewarded negatively for each collision with the prey ("every time they are eaten"). Using such reward functions, the prey group is expected to learn defense strategies and the predator is expected to learn attack strategies.

Simulation Setup

Action Space



State

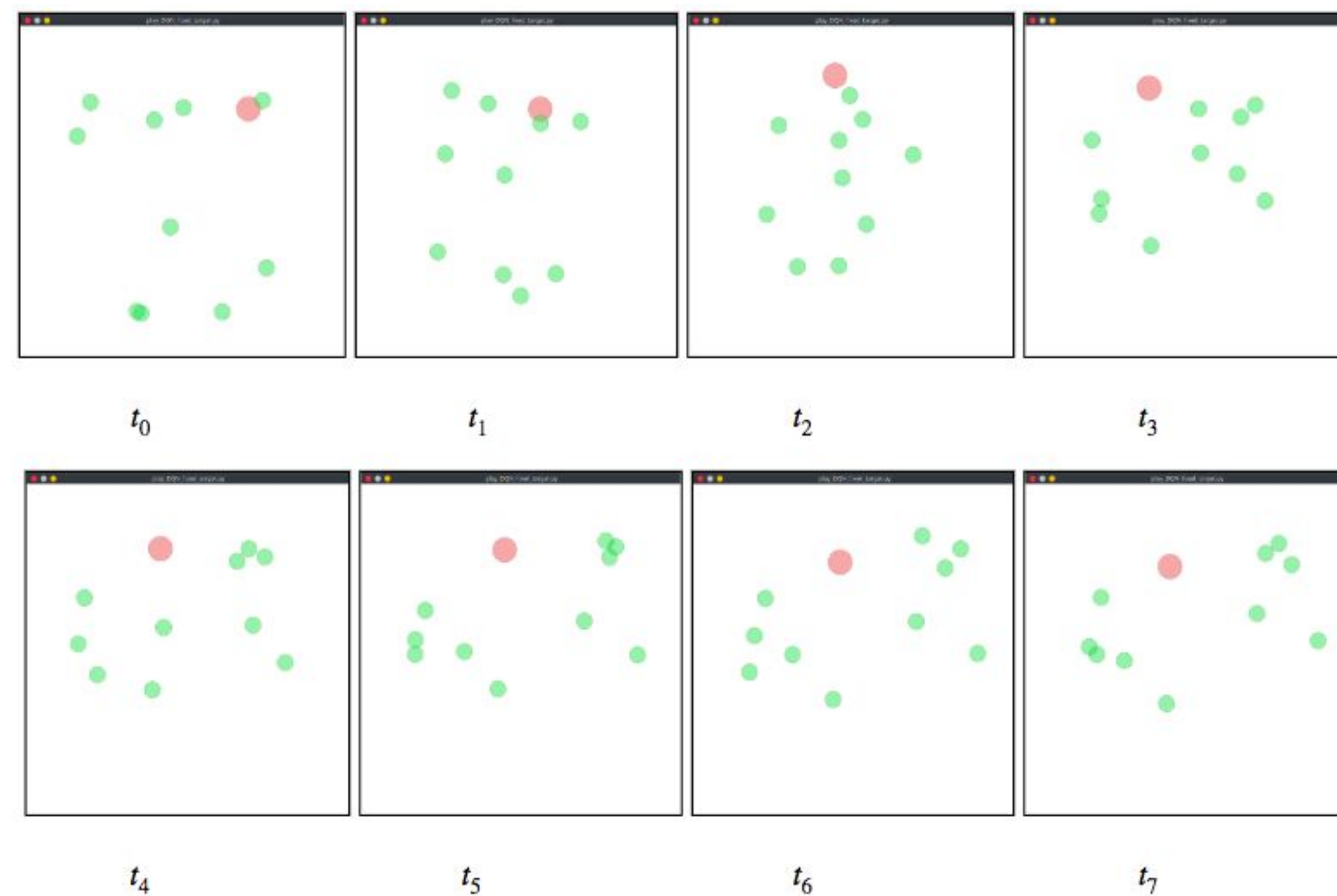
[state_i] for agent i is defined as:

State_i: [position of agent i, velocity of agent i, relative position of other agents, relative velocity of other agents]

Reward

Predator: +10 for each collision & Boundary Restriction

Prey: -10 for each collision & Boundary Restriction



Methods

Algorithm 1: Q-learning with experience replay

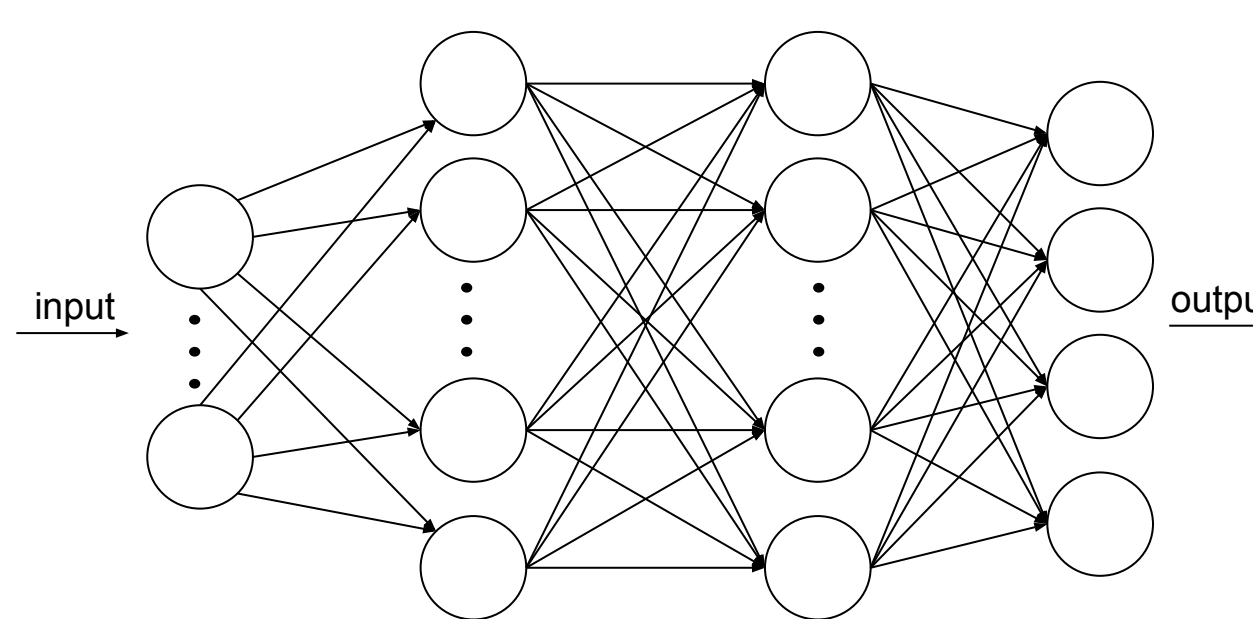
```
Initialize a replay memory to capacity N
Initialize action-value function Q with random weights θ
For episode 1:M do:
  Initialize environment to get initial state s1
  For step 1:T do:
    With probability ε select random action at
    Otherwise select at = argmaxa Q(st, at; θ)
    Execute action at and observe reward rt+1, st+1
    Store experience st, at, rt+1, st+1 in replay memory
    Set st+1 ← st
    Sample a batch sj, aj, rj+1, sj+1 from replay memory
    Set target Yj = rj + γ maxa' Q(sj+1, a'; θ)
    Perform gradient descent step on (Yj - Q(sj, aj; θ))2
  End For
End For
```

Algorithm 2: Q-learning with experience replay & fixed Q-target

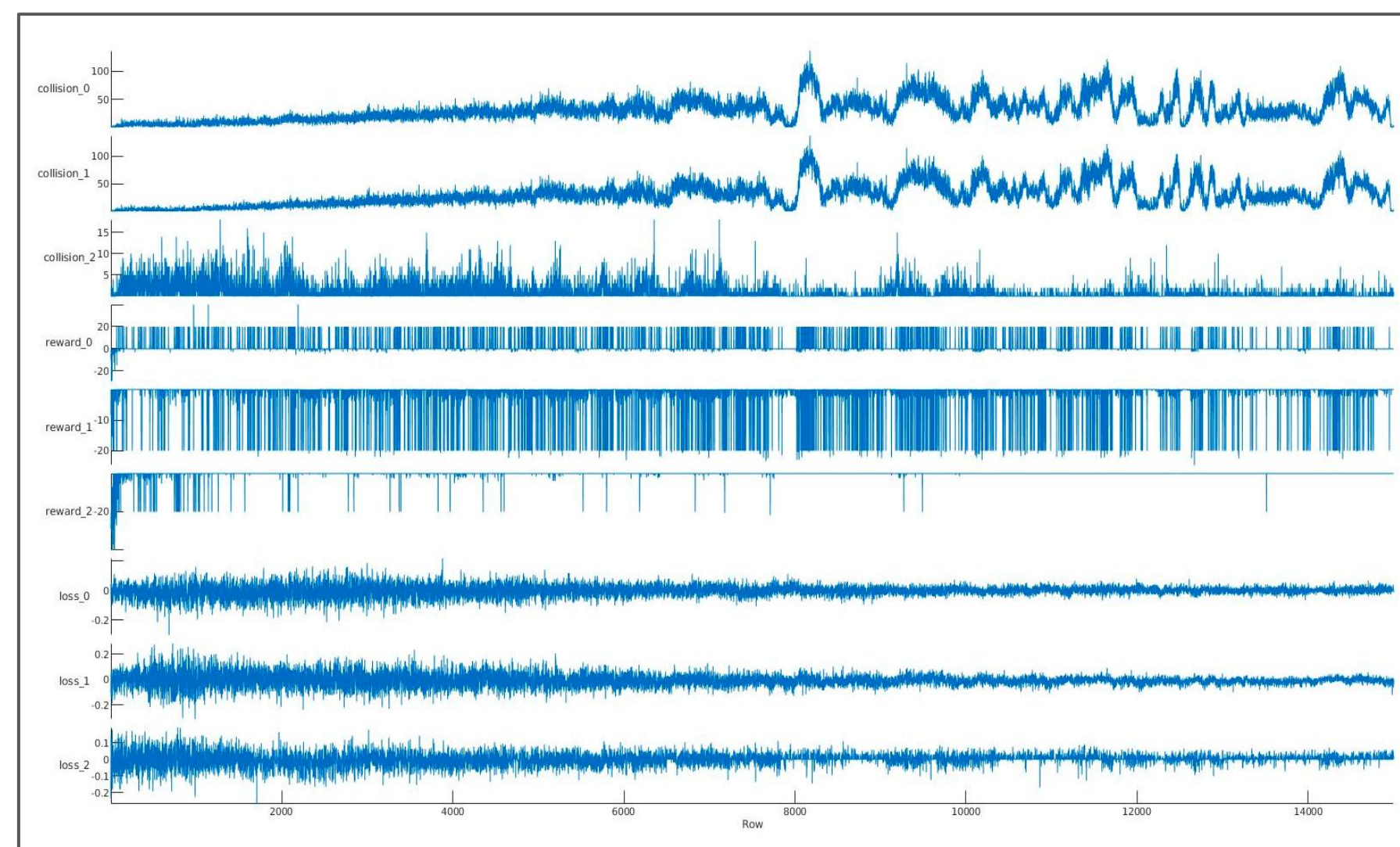
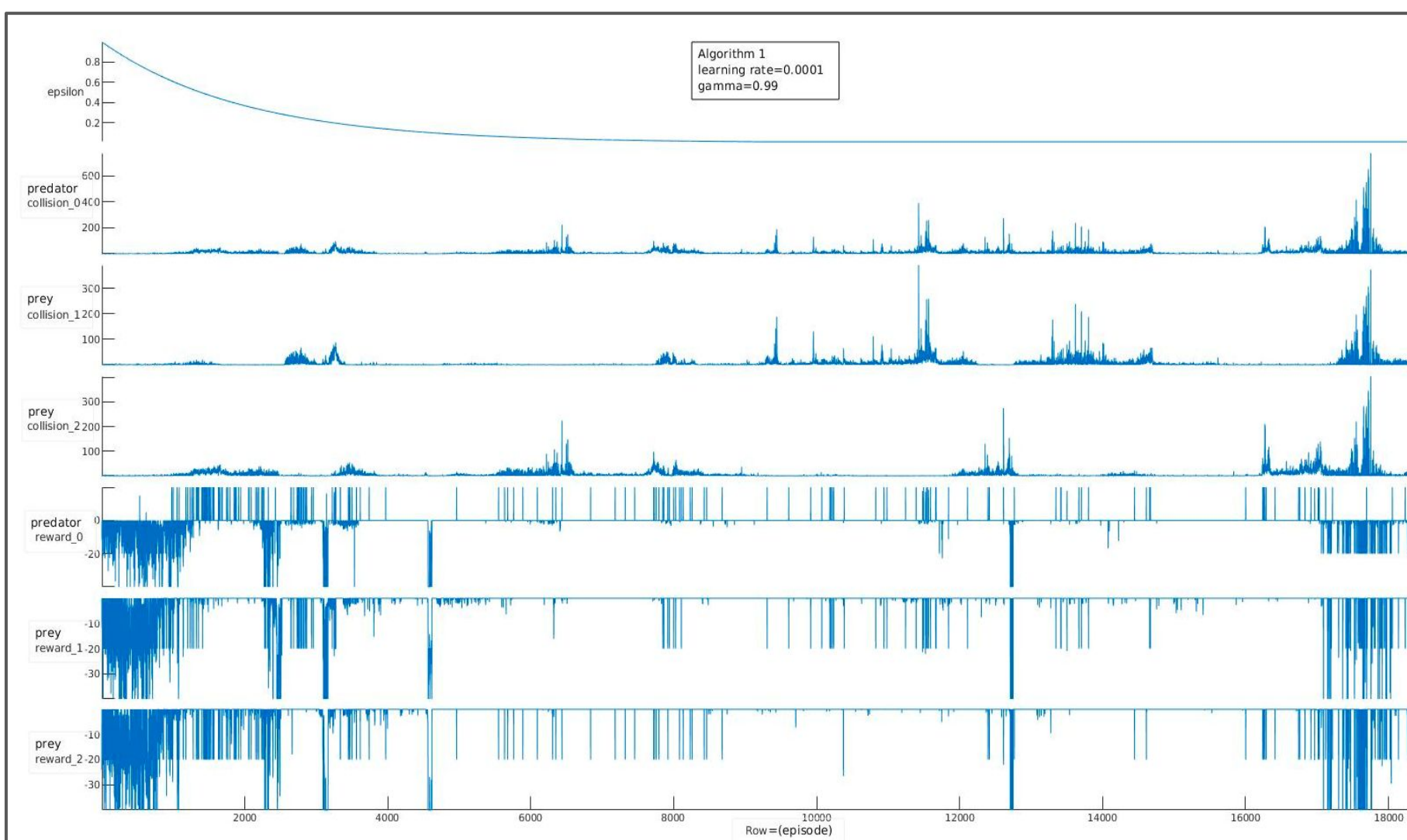
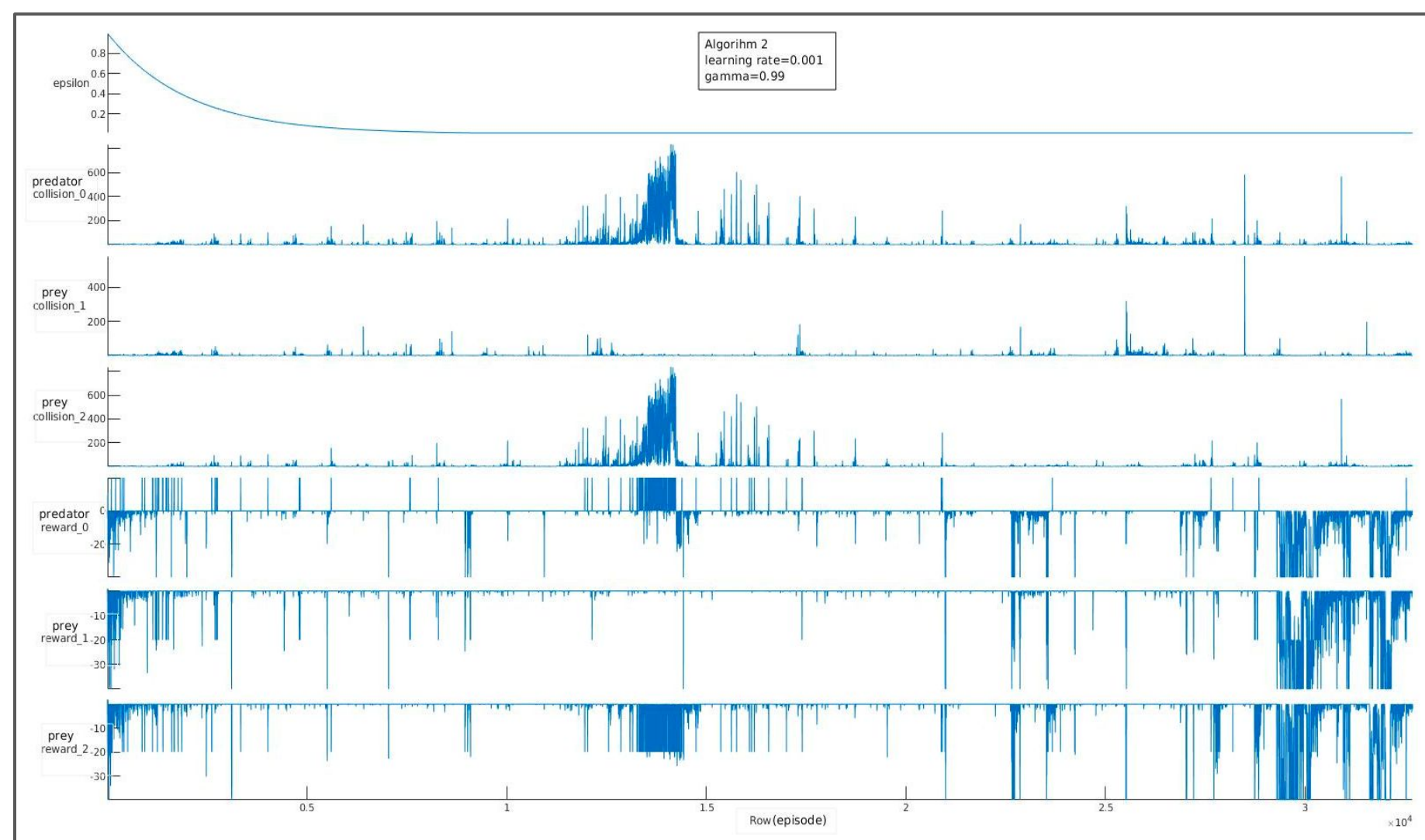
```
Initialize a replay memory to capacity N
Initialize action-value function Q with random weights θ
Initialize target action-value function Q̂ with random weights θ̂
For episode 1:M do:
  Initialize environment to get initial state s1
  For step 1:T do:
    With probability ε select random action at
    Otherwise select at = argmaxa Q(st, at; θ)
    Execute action at and observe reward rt+1, st+1
    Store experience st, at, rt+1, st+1 in replay memory
    Set st+1 ← st
    Sample a batch sj, aj, rj+1, sj+1 from replay memory
    Set target Yj = rj + γ maxa' Q̂(sj+1, a'; θ̂)
    Perform gradient descent step on (Yj - Q(sj, aj; θ))2
    Every C step update θ̂ = θ
  End For
End For
```

Algorithm 3: Monte-Carlo Policy Gradient

```
Initialize policy parametrization with random weights θ
For episode 1:M do:
  Generate an episode following πθ
  For step 1:T do:
    Set G ← ∑k=t+1T γk-t-1 Rk
    Set θ ← θ + α γt G ∇ ln π(At | St, θt)
  End For
End For
```



Results



Conclusion

In conclusion, we note that using reinforcement learning algorithms and simple reward functions for agents, it is possible for the prey agents to discover cooperative strategies that prevent the predator from "eating" them. It is also interesting to note that though the prey agents are the same in the beginning, as training proceeds, the same prey agents learn different strategies. While some are able to successfully evade predators, others are not. Such a concept is similar to the concept of "fitness" observed in real biological systems, where some individuals are able to survive better than others.