

Distributed MARL Predator-Prey Final Report

Luca Fabri

`luca.fabri@studio.unibo.it`

November 29, 2024

Reinforcement Learning (RL) is a sub-area of Machine Learning which focuses on building strategies to make decisions that maximize the future expected reward. The RL finds applications in problems concerning self-driving cars, industrial automation, and finance but also in systems that involve the interaction of multiple agents in a shared environment, where it's more specifically called Multi-Agent RL (MARL). In computational biology, it's useful to study the population of one or more intelligent species that interact with each other to build population models: in this project, a predator-prey ecosystem is implemented, using a MARL approach in a mixed environment, i.e. cooperative and competitive, where agents of the same species cooperatively make decisions to maximize their total expected reward. To achieve this goal, the MADDPG algorithm is exploited. The distribution of the system is realized by parallelizing the environments and by introducing a distributed training technique, inspired by the Mava [1] framework.

Contents

1	Goals/requirements	3
1.1	Use Cases	4
2	Requirements analysis	5
2.1	Functional requirements	5
2.2	Non-functional requirements	6
3	Architecture Design	7
4	Microservices design and implementation	9
4.1	Predator-Prey Service	9
4.1.1	Design	9
4.1.2	Implementation	13
4.2	Learner Service	15
4.2.1	Design	16
4.2.2	Implementation	17
4.3	Replay Buffer Service	18
4.3.1	Design and Implementation	19
4.3.2	Self Assessment/Validation	19
5	Deployment	20
6	Continuous Integration	21
7	Experiments	22
8	Conclusions	25

1 Goals/requirements

This project aims to implement a distributed application for training a Multi-Agent Reinforcement Learning (MARL) system. Agents' experiences are collected from multiple parallel environments and integrated into a centralized dataset, ultimately used for the learning process.

The goals of the project include:

- Development of a MARL Environment, specifically a Predator-Prey Environment, where predators are trained to catch prey, while prey to run away from predators;
- Development of a Replay Buffer, whose API allows the Environments to store the experiences of each of its agents;
- Development of a Learner Service, whose goal is to train a separate MADDPG model for predators and preys, and update the agents' policy inside each Environment.

Term	Definition
System	The <i>Distributed MARL Predator-Prey</i> application
Agent	Decision-making entity interacting with an environment through observations, rewards, and actions
Environment	General term referring to a Multi-Agent Cooperative-Competitive Environment. It is a collection of Agents inside a space, with equal or conflicting reward structure
Learner	Entity capable of training a Multi-Agent Environment through Reinforcement Learning
Replay buffer	Dataset of agent experiences
MADDPG	Multi-Agent Deep Deterministic Policy Gradient algorithm [3]

Table 1: Term's glossary

1.1 Use Cases

A user that interacts with this application is in front of three choices:

- Train the MARL system, by specifying the number of Predator-Prey Environments to run in parallel and setting a configuration file to customize the training process;
- Start a simulation. This option is only possible if a previous training phase has been carried out: each Environment will load the latest trained model and let it be used by its agents;
- Visualize the training results of the last training process, comprising a plot with neural network loss over time;
- Visualize agents of an environment in a scatterplot animation, by specifying the index of the environment.

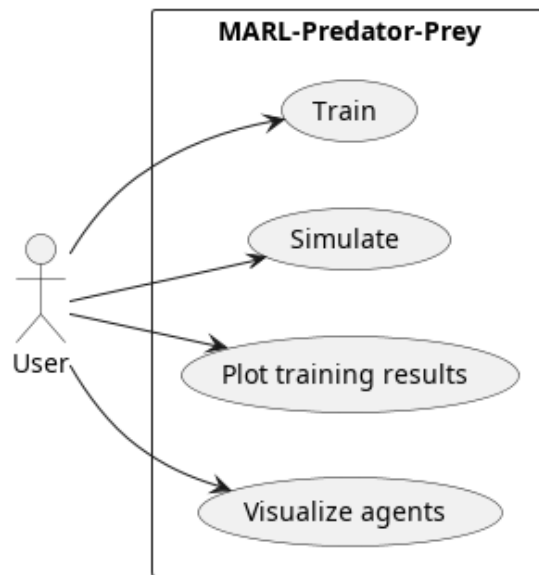


Figure 1: Use cases UML

2 Requirements analysis

In this section, the functional and non-functional requirements of the project are listed.

2.1 Functional requirements

1. System

- a) The System must accept a different set of parameters for different runs, in order to configure the training/simulation process;
- b) The System must be able to run in Train mode;
- c) The System must be able to run in Simulation mode.

2. Environment

- a) The Multi-Agent Reinforcement Learning Environment must be a mixed Cooperative-Competitive Environment with continuous state/action space, where agents are free to move in a 2D Euclidian space of w x h dimension;
- b) The Environment must be a Predator-Prey Environment, specifically.
 - i. Predators must have a reward function that increases as they approach Preys;
 - ii. Preys must have a reward function that decreases as they approach Predators;

3. Learner

- a) The Learner must update the policies of the agents inside each Environment, in a distributed way;
- b) The Learner must implement the MADDPG algorithm as the learning algorithm.

4. Replay Buffer

- a) The Replay Buffer must collect the agent's experiences of each Environment, in a centralized way;

2.2 Non-functional requirements

1. Availability

- a) Fault tolerance: The system should be able to recover from failures and continue to operate.

2. Deployability

- a) Portability: The system should be able to run on different platforms.

3. Modifiability

- a) Extensibility: The system should be able to add new features easily.
- b) Maintainability: The system should be easy to maintain and update.

3 Architecture Design

The **microservices architecture** has been chosen for the system. This type of architecture consists of decomposing the project into smaller parts so that each one can be deployed independently. This will make it easier to add new features to the system, scale the Environment, and generally maintain the project.

The system is decomposed into three microservices:

1. **Predator-Prey service.** This component is responsible for managing a single Environment. Supports both training mode and simulation mode;
2. **Learner service.** This component is responsible for training the System;
3. **Replay Buffer service.** This component is responsible for integrating all agent experiences into one single dataset.

For communication, the Replay Buffer service will expose a REST API for batching and storing data. The Environment-Learner link instead uses the Publish/-Subscribe pattern, to send/receive data, made possible through a message broker.

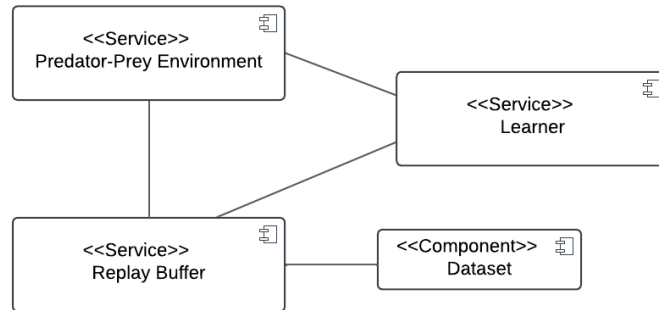


Figure 2: System's Architecture

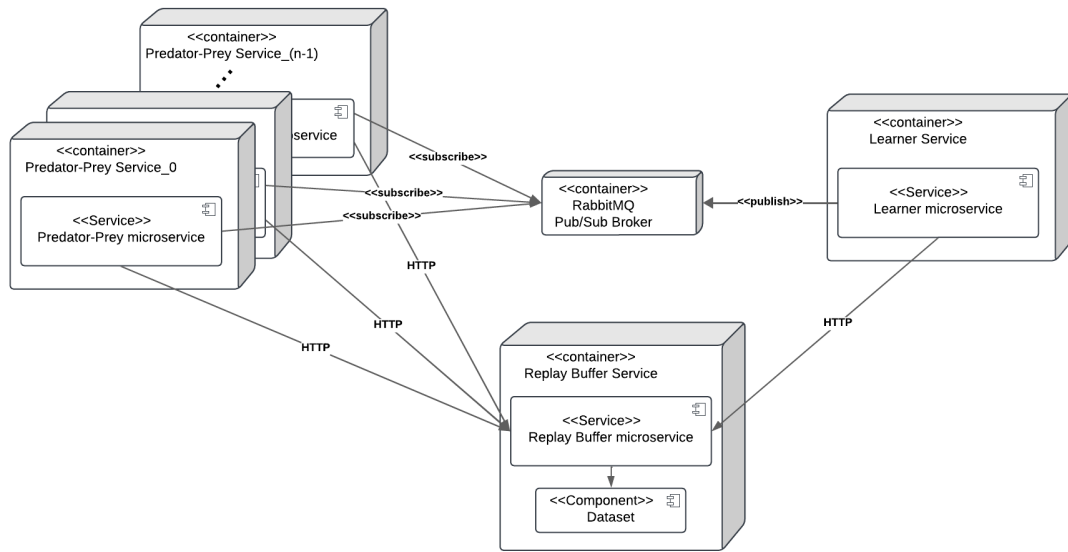


Figure 3: Training mode Deployment View

4 Microservices design and implementation

4.1 Predator-Prey Service

The Predator-Prey Service is responsible for managing the Predator-Prey environment. It allows to start the environment in two modes: Train mode or Simulation mode.

1. In the first mode (Train), it interacts with a distributed Replay Buffer to store agents' experiences and subscribes to policy updates coming from the Learner Service;
2. In Simulation mode, it uses the latest policy (saved from a previous training phase) to start a simulation. In this case, it accepts a random seed as input to control the initial positions of the agents.

4.1.1 Design

The service follows the MVC pattern to better separate the code responsibilities:

- The Model consists of the `Environment` and `Agent` classes. In particular, the `Environment` models a 2D Euclidian space using two integers, one for specifying the x -dimension and another for the y 's. It's composed of a list of `Agents`, each one with its coordinate within the space and velocity (both v_x and v_y components);
- The Controller contains the `EnvironmentController` and `AgentController` interfaces and implementations. The first one collects single agents tuples (State, Action, Reward, Next state) and stores the joint tuple inside the Replay Buffer, while the latter is responsible for managing the single `Agent`. It implements the reward and done functions and contains an `AgentPolicyController`.
- The View only contains the service's entry point.

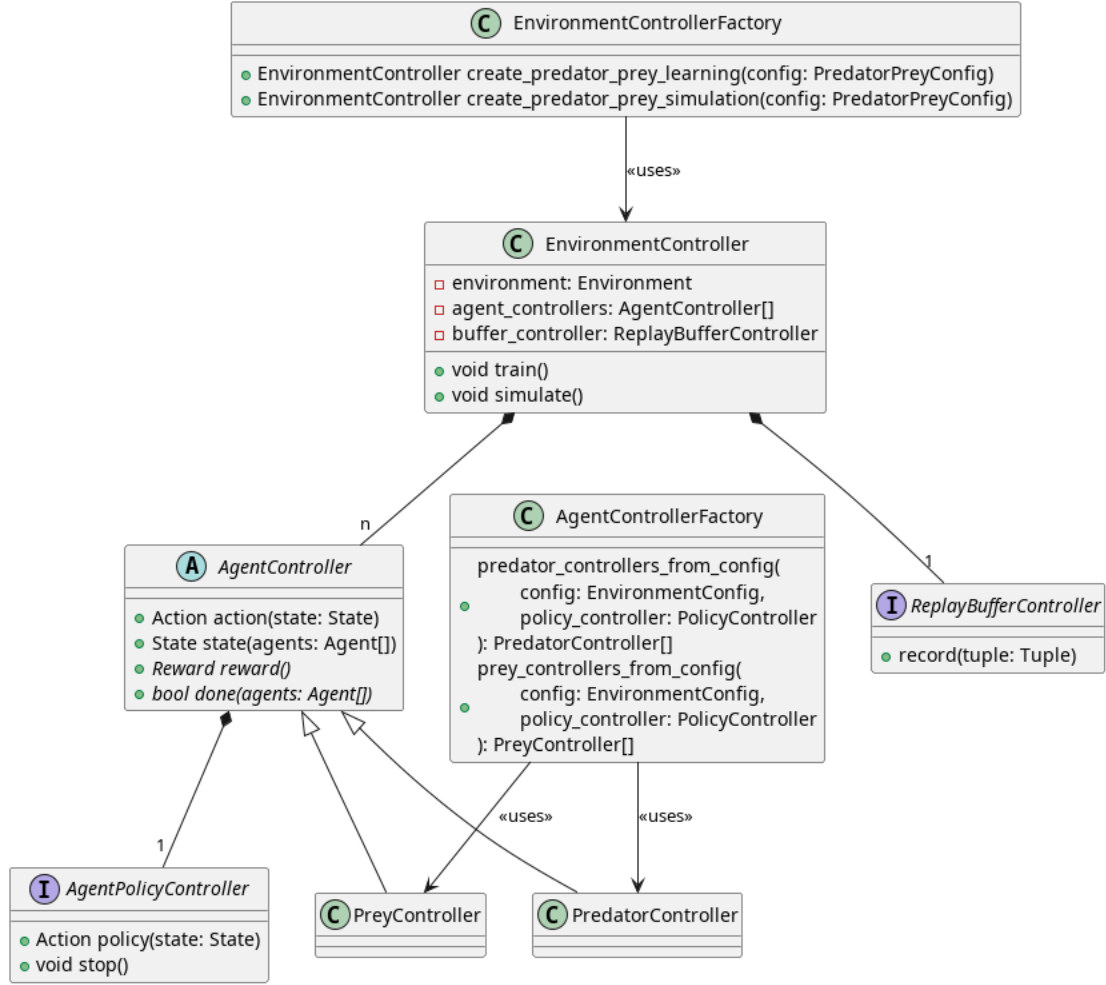


Figure 4: Controllers

As previously stated, **AgentController** manages the single **Agent**. Whether it's a predator or a prey, the state/action space is the same:

- The state reflects an agent's observation. It's composed of a set of numbers representing the agent's view of its surrounding area. More precisely, the observation is modeled as a pencil of line segments passing through the agent's coordinate. The line segments have a specific length that represents the visual depth of the agent.

Using this kind of observation space allows us to get the same number of values for different agent steps. This is necessary for the neural network that will train the system, as it requires the input to always have the same shape. That being said, suppose that agent i has position (x_i, y_i) the pencil of lines is defined by the following equation:

$$(x - x_i) * \sin \alpha = (y - y_i) * \cos \alpha \quad (1)$$

where $\alpha \in [0, \pi]$. The set of α parameters is chosen to be a set of evenly spaced samples so that the lines form the same angle.

For each line segment, a single observation value in $[0, 1]$ is produced, which is:

- the distance to the intersection point if another agent of a different type (preys for predators, predator for preys) intersects the line segment, normalized by the visual depth;
- set to 1 if no agent of a different type intersects the line segment.

Actually, each line segment is divided into two parts of the same length, using the agent's coordinate as the starting point. This allows the agent to discriminate if the intersected agent is for example in front of him, or back to him. Indeed, if we were using a single value for each line segment, the agent wouldn't be able to know the direction of the one observed.

In the following figure, the agent's observation is shown as an exemplification. For better readability, the pencil of line is shown only for agent a_0 .

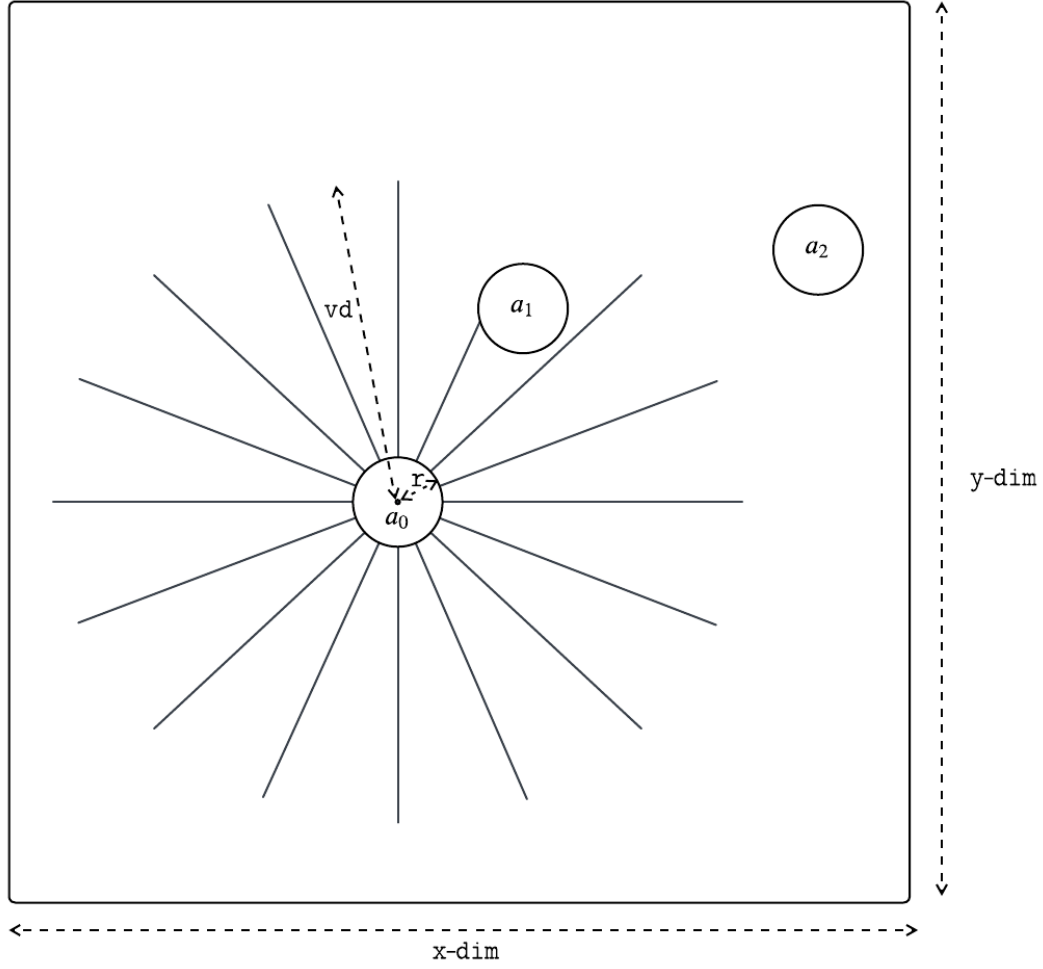


Figure 5: Agent's Observation

Both visual depth vd , radius r , x_dim and y_dim are configurable by setting their values in a configuration file, as well as the number of line segments `num_states`.

- The action instead is composed of two numbers, one representing the speed $\|\vec{v}\|_2$ and the other the angle θ . Given the agent- i 's coordinate (x_i, y_i) and time delta Δt , it's easy to compute the next agent's position.

While the space/action space is the same in both predators and preys,

`PreyController` and `PredatorController` have different reward functions. Both of them are linear and produce a value in $[-1000, 0]$, computed based on the closest agent (of the other type) inside the observation: the predator one increases as the closest prey gets closer, the agent one increases as the closest predator moves away.

Note that if no agent appears inside the observation space, predators will have the minimum reward (-1000), while the preys the maximum (0).

4.1.2 Implementation

The controller implementations allow the service to communicate with the others that compose the distributed application.

Specifically:

- The `RemoteReplayBufferController` implements `ReplayBufferController` and allows to record the joint tuple to the distributed Replay Buffer. It leverages the Requests Python library to send HTTP Post messages to Replay Buffer Service;
- The `AgentPolicyController` instead has two implementations, one for the Train mode and one for Simulation mode. The first one contains an `ActorReceiverController`, that subscribes for policy updates coming from the Learner Service, while the other loads an existing policy from a file, saved from a previous training phase.

For receiving an agent's policy the implementation uses Pika, a RabbitMQ client library for Python, by subscribing for both Predators and Prey policy updates.

The agent's observation inside `AgentController` is instead computed by leveraging the Python `z3` library, a powerful tool for SMT solving. Since solving non-linear equations with this tool is too time-intensive, the agent shape has been changed to a square (of side length $2r$) to speed up the computation.

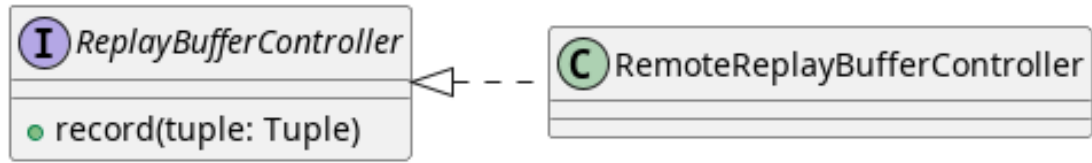


Figure 6: Buffer

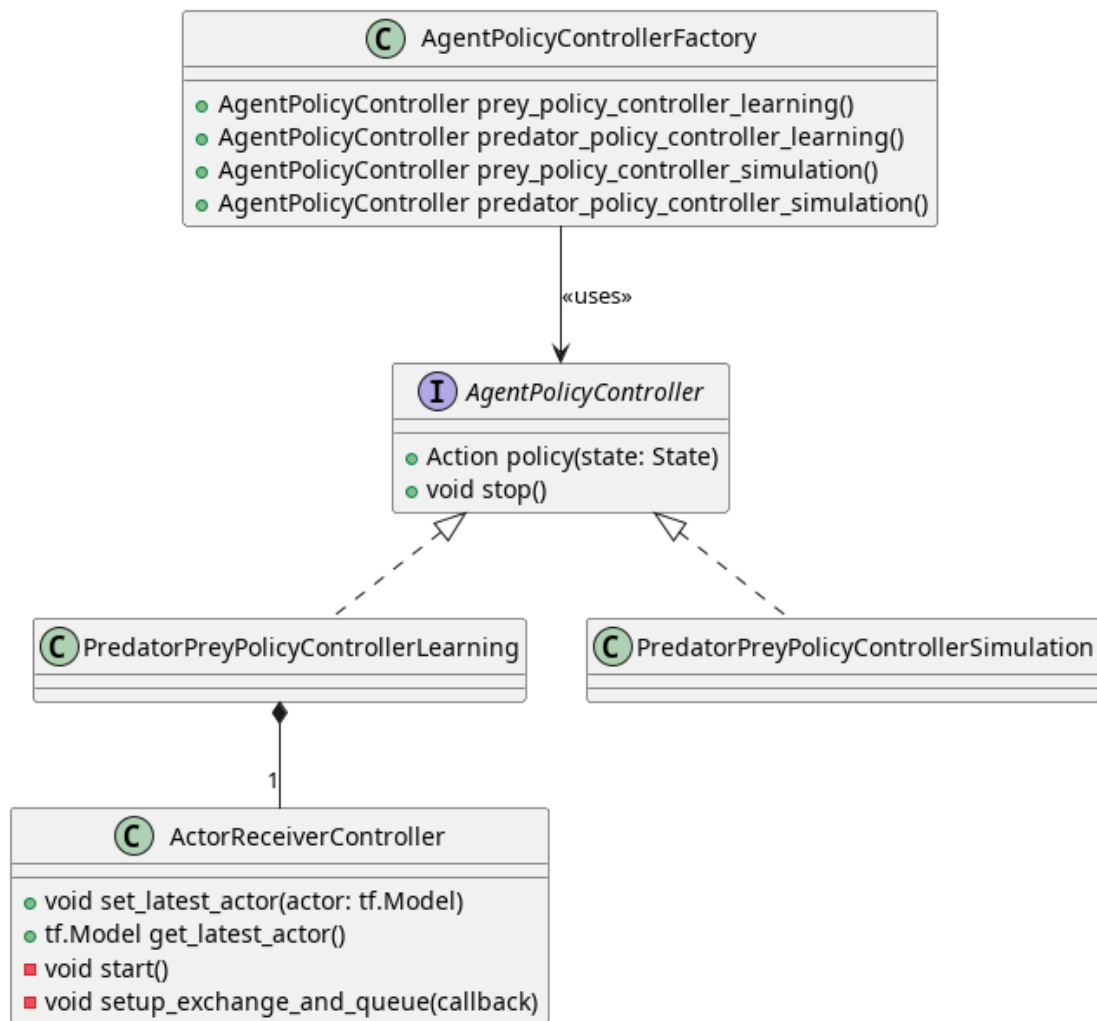


Figure 7: Policy

4.2 Learner Service

The Learner Service is responsible for training the agents belonging to the different MARL environments.

From the Learner’s perspective, the environment is unique, as it samples data from the centralized Replay buffer where the agent’s experiences are integrated into the same database.

Specifically, the Learner is based on the MADDPG [3] training algorithm, the multi-agent version of DDPG [2], which in turn uses an Actor-Critic model. The Actor-Critic structure decomposes the learning procedure into two parts:

- The Actor, given an agent’s observation, decides which action should be taken, so it learns the so-called policy;
- The Critic, also known as the value function, evaluates how good the chosen action is to guide the Actor towards decisions that lead to higher expected rewards.

This algorithm relies on the **centralized learning** and **decentralized execution** framework:

- During training, all agents have a centralized Critic network that takes in input the joint observation and the joint actions of all agents and produces the Q -value for the respective agent. Specifically, the Q -value represents the expected cumulative reward for the agent taking that particular action given the state. This value is used in turn to train the Actor network;
- During execution, each agent can rely on the latest trained Actor model, as it is a neural network that takes an observation and produces an action. In this phase, the Critic network is not needed anymore.

That being said, in order to let the agent operate, the Learner only needs to send the Actor model over the network to each distributed Predator-Prey environment.

This decomposition allows each environment to not depend on the Learner in Simulation mode, as each agent can move thanks to the latest Actor model received during the Training phase.

In the following figure, a simplified schema of the framework is shown. Rewards r and Next states s' are not shown inside the Learner, but they are used for minimizing the loss of the Critic network, achieved through a form of temporal difference learning (TD-learning).

Another sketch of the framework, extracted from the original MADDPG publication can be found [here](#). For more insights about the algorithm, please refer to the original paper or navigate to this [link](#) for the pseudocode.

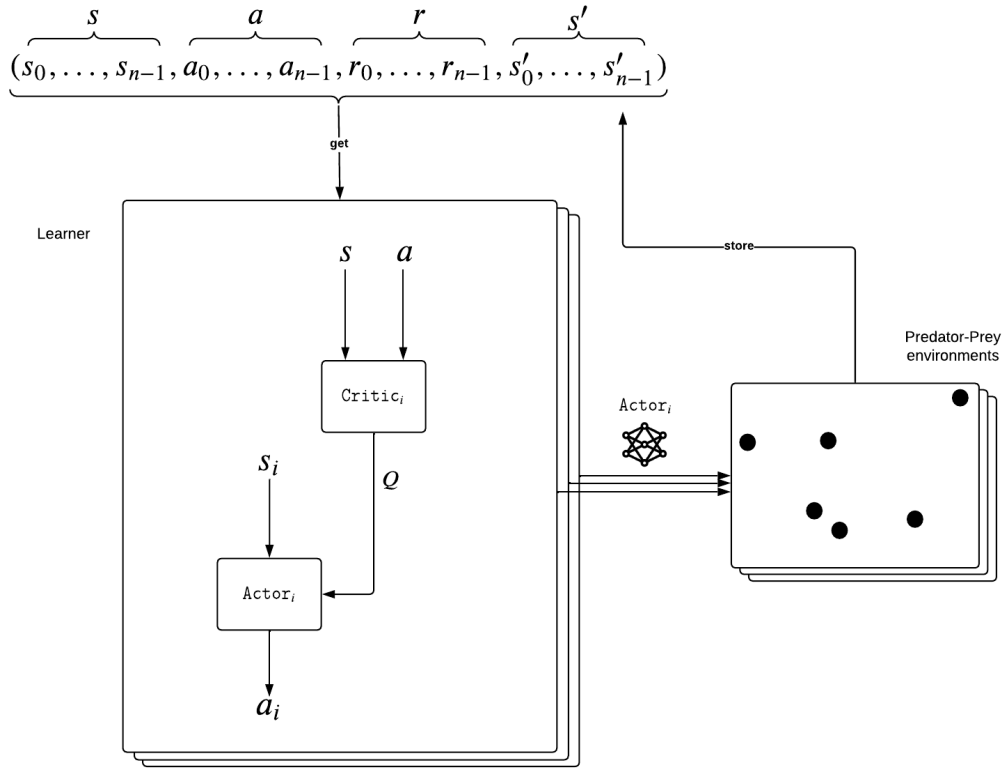


Figure 8: Learner Actor-Critic schema

4.2.1 Design

The design of Learner Service follows the MVC pattern:

- Model contains the **Actor** and **Critic** network models and the configuration parameters;
- Controller comprises three controllers: **ReplayBufferController**, used for batch data from the Replay Buffer, **ActorSenderController** for sending the Actor networks and **LearnerController**, whose implementation is responsible for training the system.
- View only consists of the service's entry point.

The neural networks structure is the following:

- Critic:
 - Input:
 1. Array of joint states s , of shape $(\text{num_states} * \text{num_agents}, n)$, where n is the batch size;
 2. Array of joint action a , of shape $(\text{num_actions} * \text{num_agents}, n)$, where n is the batch size;
 - Inner layers: three fully connected layers with ReLU activation function of 256, 128, 128 neurons, respectively;
 - Output: a real number representing the Q -value.
- Actor:
 - Input: Array of shape $(\text{num_states}, n)$ where num_states is the number of line segments used for the observation and n batch size;
 - Inner layers: two fully connected layers with ReLU activation function of 128, 64 neurons, respectively;
 - Output: two real numbers representing the speed $\|\vec{v}\|_2$ and the angle θ , limited in $[-1, 1]$ by using tanh activation function.

4.2.2 Implementation

The controller's implementations are:

- **RemoteReplayBufferController** to batch data from the distributed Replay Buffer Service;

- `PubSubActorSenderController` to update the agents' policies using the Publish/Subscribe pattern. As for the Predator-Prey Service, it leverages the Pika library to publish messages;
- `MADDPGLearnerController` that implements the MADDPG algorithm, using Tensorflow library. The implementation is based on the project MADDPG-Keras.

The `Actor` and `Critic` models inside the `Model` instead use the Keras API.

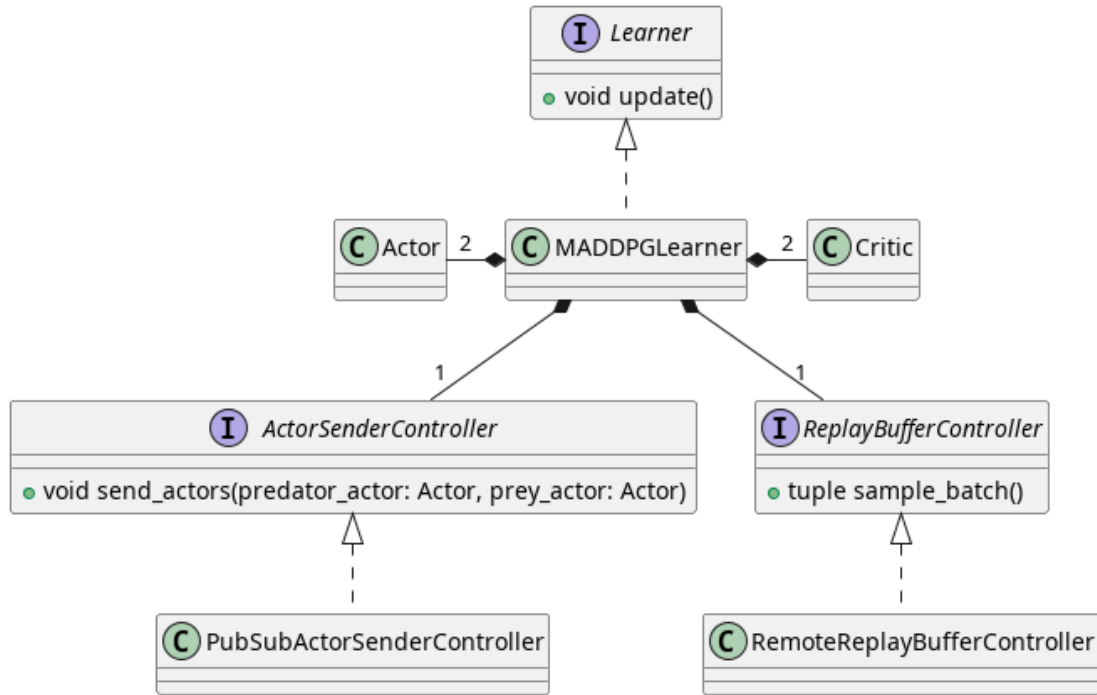


Figure 9: Learner class UML

4.3 Replay Buffer Service

The Replay Buffer Service is responsible for storing the agents' experiences inside a dataset, used for training the System.

It provides a simple and intuitive API to get and store data.

4.3.1 Design and Implementation

Replay Buffer Service's design is minimal. It comprises a `ReplayBufferService` class that manages the exposed API and a `DataBatchValidator` that verifies the received data batch has the correct structure.

In particular:

- The exposed API is an HTTP ReST API that allows to extract or record a data batch:
 - A HTTP POST request to `/record_data/` containing a data batch as Json records data to the replay buffer;
 - A HTTP GET request to `/batch_data/<size>`, where size is the number of rows to batch is used to extract data.
- `DataBatchValidator` verifies that the received data batch to store has a compliant structure. Each row should contain:
 1. The joint state: the observation of each agent;
 2. The joint action: the action of each agent;
 3. The joint reward: the reward of each agent after executing the respective action;
 4. The joint next state: the observation of each agent after executing the respective action.

Replay Buffer Service is implemented in Python3 using Flask for the HTTP API with Gunicorn for production deployment. The OpenAPI specification is available [here](#).

4.3.2 Self Assessment/Validation

Replay Buffer Service contains tests to verify the HTTP API works as expected. They are implemented using Behave, a Python library for Behavior-Driven Development (BDD).

The scenarios test both `record_data` and `batch_data` APIs and are available at the following link.

5 Deployment

The System is deployable both in Train mode and Simulation mode, through a Docker Compose file that runs all microservices at once.

To do so, navigate to the Bootstrap repository and follow the provided instructions.

In Train mode, all the System's microservices are deployed with the following dependencies.

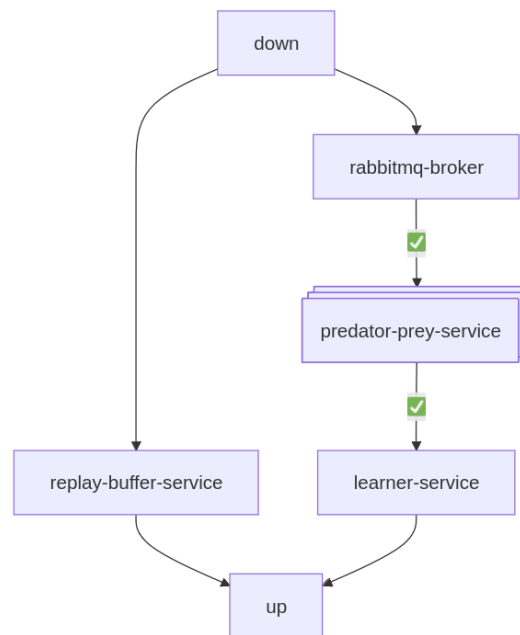


Figure 10: Train mode - Microservices deploy

In the `docker-compose` file, it's guaranteed they are deployed in this sequence by leveraging healthchecks and `depends_on` attributes.

In Simulation mode, only multiple `predator-prey-service(s)` are deployed. They are run in parallel with no dependency on each other.

6 Continuous Integration

In all the microservices repositories a Continuous Integration pipeline is built. The pipeline ensures the dependencies installation is always successful and that the code is properly formatted.

The same CI/CD pipeline is adopted by all the microservices, with an exception for `replay-buffer-service` where tests are introduced in the `build` job to ensure they always pass. The CI workflows are the following:

1. `ci.yaml`. It is in turn composed of the following jobs:
 - a) `build`: is responsible for checking if the project dependencies are successfully installed. A matrix of OS is provided: Ubuntu, MacOS, Windows;
 - b) `format`: checks if the code format is correct leveraging Ruff formatter.
2. `deploy-image.yaml`: builds a Docker image of the microservice and publishes it to GitHub Packages;
3. `gh-pages.yaml`: builds the Code documentation of the microservice using Sphinx and publishes it to GitHub Documentation;
4. `release.yaml`: Produces a release when a Git tag is produced and publishes the code to GitHub Releases.

7 Experiments

The System has been trained for a total of 7h 12m 3s, in a Manjaro Linux 24.1 OS, Kernel version 6.6.54-2, using a CPU Intel(R) Core(TM) i7-8550U @ 1.80GHz and 16GB RAM.

The Critic network loss plot has been generated and shown in the following figure and a Predator-Prey-Service-0 animation, comprising the latest 500 steps is available at the following link.

By visually analyzing the Predator-Prey animation and the Critic loss over time, it can be concluded that the algorithm doesn't converge. Both Predators and Preys Critic loss suddenly increases at nearly iteration 3000th, and slowly decreases after that point. However, after 7h the loss is still nonzero. Moreover, the animation doesn't show any agent behavior: the predators seem to not chase preys, preys seems to not run away from predators.

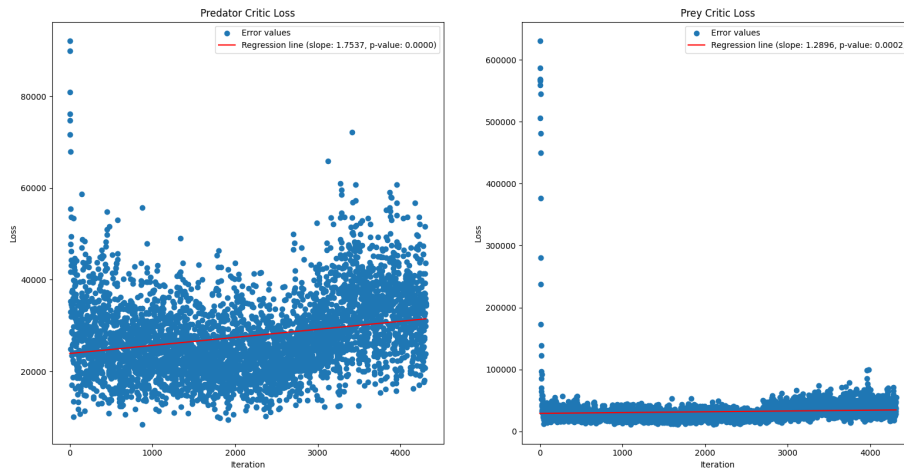


Figure 11: Predator and Prey Critic loss

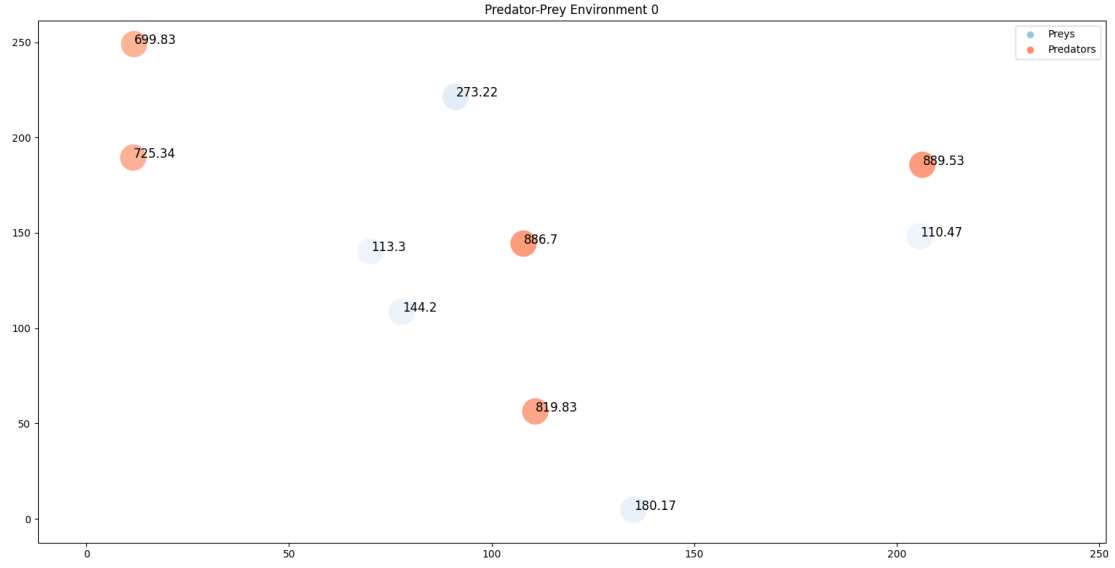


Figure 12: Animation frame - Numbers near agents represent their last reward.

The non-convergence of the system may result from multiple factors:

1. Rate of collected data rows per time unit.
 - The agent’s observation function, of both predators and preys is relatively complex: using the hardware setup listed before, each environment step takes nearly ~ 30 seconds. For ~ 7 h of training and 5 environments in parallel, only 4,2k rows are collected inside the Replay Buffer. This data may be insufficient for the Learner to learn a definitive policy: to collect sufficient data the system may be trained for more time.
 - MADDPG implementation is based on the existing project MADDPG-Keras. The author states: “It takes around 20 hours to train 3 agents in 2 pursuer-1 evader environment for 3000 episodes (100 steps in each episode) on single i5-113G7 processor”.
- Due to a lack of time and hardware resources, the student couldn’t run the system for the time needed to collect 300k rows.

2. Algorithm implementation and parameters. Modifications had to be applied to the MADDPG project linked before. It cannot be excluded that the code has some implementation errors. The algorithm is difficult to test, so errors may have been introduced. Moreover, parameters like the discount factor and the learning rate have been set to those proposed in the DDPG Pendulum example provided in the Keras website. These parameters should be further verified to be appropriate for this kind of environment.

8 Conclusions

Future experiments should be conducted by training the system in a dedicated machine and for more time. Indeed, as stated in the previous section, I had the impossibility to run the application for further time.

Help from domain experts is also well appreciated to verify the code is compliant with MADDPG algorithm. Implementing MADDPG has been proven challenging due to the lack of available open-source projects for comparison.

References

- [1] Ruan de Kock et al. *Mava: a research library for distributed multi-agent reinforcement learning in JAX*. 2023. arXiv: 2107.01460 [cs.LG]. URL: <https://arxiv.org/abs/2107.01460>.
- [2] Timothy P. Lillicrap et al. *Continuous control with deep reinforcement learning*. 2019. arXiv: 1509.02971 [cs.LG]. URL: <https://arxiv.org/abs/1509.02971>.
- [3] Ryan Lowe et al. *Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments*. 2020. arXiv: 1706.02275 [cs.LG]. URL: <https://arxiv.org/abs/1706.02275>.