

Reinforcement Learning Project

Neural MMO by DQN, PPO, and Curriculum

Xianjin GONG, Hangao LIANG, Gaoyuan ZHOU
(Github repository: <https://github.com/jimmily98/INF581-project>)

Abstract—So far, we have gained experience of basic reinforcement learning methods including Q-learning, SARSA, Deep Q-learning (DQN), and Policy Gradient (PG) methods. All of them were only applied on single agent environments with small observation and action space such as Cartpole or Lunarlander provided by gym. In this project, we deploy DQN to train multiple agents in a Massively Multiplayer Online-Game (MMO) environment. The result is compared with the one by Proximal Policy Optimization (PPO). And we will also examine how curriculum can be designed to influence the learning process.

I. INTRODUCTION

In this project, we are interested in the performance of DQN model in an MMO game environment with multiple agents.

The chosen MMO environment comes from the competition hosted by NeuralIPS. It is close for submission now. But because this environment is well constructed following the frame of PettingZoo, it is easy for us to adapt our existing code. And because it provides tracks for both reinforcement learning model construction and curriculum building, it fits our interests quite well.

Compared with tabular methods, DQN is good at expressing policy with large observation and action space by using neural networks. It can directly process visual information, such as the image input from Atari games, by making use of convolutional neural networks (CNNs). It can understand temporal information by stacking several frames as one input, such as learning the movement of the ball in the game of Pong.

However, so far, we only implemented DQN with single agent environment. If we have several agents interacting with each other in one environment, especially for the case that they need cooperate, how should we train them with DQN?

One way is to train a central brain which gives orders to each agent. But in this method, the space dimension for the central brain increases a lot and we did not find a good solution to solve that.

Another intuition from our group is that we can train each of them separately for every time-step. And when we train one of them, we regard the rest agents' information as part of the environment. Every agent uses the same policy. We make them cooperate by tweaking the reward function.

During experimenting DQNs, we suffered a lot from their instability. This instability partly comes from the fact we were bootstrapping and optimizing toward a moving target, which was relieved by using double networks. But it also comes from the fact that DQN is an online method. Although batch updating can decrease certain noise during optimization,

higher batch size will increase computation time for each time-step, which results in a new trade-off between accuracy and cost.

After researching, we find that PPO is more stable than DQN and it is currently the state-of-the-art (at least the foundation of state-of-the-art) for reinforcement learning. As a policy-based method, PPO directly outputs the probability distribution of actions given the current state. In this way, it handles continuous action space better than DQN which is a value-based method accompanied with a manual policy such as ϵ -greedy policy.

More importantly, PPO is more stable because it cares about how much change was made between the old and new policy. PPO defined a ratio between the probabilities of one action by the new policy and the old policy. It clips the gradient based on the difference between this ratio and one. If the ratio is too large or too small, the clipped gradient prevents the optimization process from stepping away from the optimal path. Thus it avoids stepping into higher gradient region, which means higher possibility to overshoot and higher variance.

There has already been one PPO implementation provided by our MMO game environment. We use this implementation and compare it with our own DQN implementation.

The last part that we are interested to dive in is the design of curriculum.

Our MMO environment is quite complicate that it involves harvesting food, building armors, trading in market, combating with NPCs, upgrading professions, and etc. To speed up the convergence process, one possible solution is to build a learning path (or curriculum) for agents to follow. We can request them to learn how to harvest food and combat against enemies first, and postpone the optimization of trading actions until a later stage. We would like to create several curriculum and see how this will affect the learning process.

In our project, we successfully implemented DQN for NeuralMMO environment. But we did not succeed to train agents well compared to PPO with our limited resources. We will discuss in our report why PPO is much better when the scale of environment is large. We also tested two curriculum. One is simple and general with only a single task to ask agents to stay alive as long as possible. Another one is comprehensive that it is composed by many intermediate small tasks such as eating food and fighting enemies with lower level. We will discuss why the second one is better than the first one.

II. BACKGROUND

A. Reinforcement Learning Model

Brief and essential details of the MMO environment is as the following (we use similar API as the Gymnasium):

- Map: 128x128 tiles. Tiles types include resource and non-resource, passable and obstacle.
- Agent visibility: 15x15 tiles centered at the agent with up to 100 entities on it.
- Observation space: Action Targets (mask indicating target of some action); Own ID; Current Time Step; Entity List in Vision; Inventory; Market; Task; Whole Map
- Action space: Move; Attack; Manage Inventory; Give; Sell; Buy; Communications

We implemented the DQN model from [3] with modifications:

ALGORITHM 1: Deep Q-Learning

```

Initialize replay memory  $D$  to capacity  $N$ ;
Initialize network list  $Q$  with random weights  $\theta$ ;
Initialize target network list  $\hat{Q}$  with weights  $\theta^- = \theta$ ;
for  $episode=1, M$  do
  for  $t=1, T$  do
    for each agent do
      for each action type do
        With probability  $\epsilon$  select a random action  $a_t$ ;
        otherwise select
           $a_t = \operatorname{argmax}_a Q[action_{type}](s_t, a; \theta)$ ;
        Execute  $a_t$ ;
      end
    end
  for each agent do
    Observe reward  $r_t$  and next state  $s_{t+1}$ ;
    Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $D$ ;
  end
  Sample random mini-batch of transitions
   $(s_j, a_j, r_j, s_{j+1})$  from  $D$ ;
  Set  $y_j =$ 
     $\begin{cases} r_j & \text{for terminal } s_{j+1} \\ r_j + \gamma \max_a \hat{Q}(s_{j+1}) & \text{for non-terminal } s_{j+1} \end{cases}$ ;
  Gradient descent on  $(y_j - Q(s_j, a_j; \theta))^2$ ;
  Every  $C$  steps reset  $\hat{Q} = Q$ ;
end
end

```

Our algorithm is different from the original one for three reasons. First, we do not use the image of the game as the observation. Therefore, we do not require a pre-processing step for each state. Second, we have 128 agents in the environment. Therefore, we have to go through each agent and store their timestep replay in the buffer. Third, for each agent, there are 9 actions with 13 action parameters to choose (action like attack requires two parameters, attack type and attack target). So, we use a list to store 13 networks for computing q-values for each action parameter.

There are two points to emphasize here. One is that we have to perform all actions of all agents before observing the reward and the next state. This is because theoretically, all agents perform the action at the same time. Another one is that because all agents share the same neural network, for each timestep, we only perform gradient descent once.

For the PPO algorithm, the implemented one by MMO environment is not different from the one introduced in [4]. The major adaptation here is that trajectories are collected from multiple agents instead of a single agent. The objective function is composed of three parts:

$$L^{CLIP+H+V} = L^{CLIP} + h^H - vL^V$$

where

$$L^{CLIP} = \hat{\mathbb{E}}_t[\min p_t(\theta)A_t, \text{clip}(p_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t]$$

$$p_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$$

$$L^V = \hat{\mathbb{E}}_t[(V_\omega(s_t) - V_t^{target})^2]$$

and the third part H is the entropy bonus whose definition is quite complicated and can be found in [5].

B. Reward and Task Management

For environments we have experienced so far, they usually have only one final optimization target and one reward function. In our case, we can define this target as staying alive as long as possible. We can encourage our agents to achieve that by giving them one point for each tick (game frame) if they are still alive. However, because our environment is too large, this reward function may produce unwanted strategies such as agents always hiding. These strategies may seem to have a good performance at the beginning, but if we allow more timesteps, we will see that they are local minimum in our optimization space.

To solve the problem, we can use curriculum composed by multiple tasks. It is a task list predefined by the environment (or can be us). Before each round (episode) of simulation, one task is randomly sampled from the list. Players receive rewards each time when they accomplish the task.

One example task is to collect 5 level-3 shard (which will be dropped by a level-3 NPC). In this round with fixed maximum timestep, every time one agent collects 5 level-3 shard, the team receive rewards. To achieve that, agent should first learn how to upgrade so that it can defeat a level-3 NPC.

Therefore, a good task list (curriculum) should teach agents different kinds of skills to survive in the environment. Trained agents will be finally assessed by some unknown task list proposed by the competition organization. The idea here is that we never know what task is contained there. So, we have to define a list which can be generalized well such that our agents can accomplish unseen tasks on unseen maps.

III. METHODOLOGY

A. Environment

To define our reinforcement learning algorithm, we first want to provide a detailed description of our environment, including its observation spaces, action spaces, and reward functions. For each agent, we can break down its observation space into eight components: Action Targets, AgentId, CurrentTick, Entity, Inventory, Market, Task, and Tile. The specific meanings of each component are outlined in [7]. We will not delve into the detailed composition of each component here for brevity.

```
observation_space(agent_id) = {
  'ActionTargets', # Action masks for the
    → action space items
  'AgentId', # Id for each agent
    → (including NPCs)
  'CurrentTick', # Current tick (0~1024)
  'Entity', #
  'Inventory', # Agent's inventory, with
    → item's type and level
  'Market', # Market for weapons and tools
  'Task', # Pre-defined task per tick
  'Tile' # Agent's location
}
```

Regarding the action space, each agent can take multiple actions per tick—one from each category. Each action accepts a list of arguments, where each argument is a discrete variable (such as a standard index for direction, or a pointer to an entity like an inventory item or agent). The action space for a single agent is divided into ten parts: Move, Attack, Target, Use, Destroy, Give, GiveGold, Sell, Buy, and Comm. The meanings of each part are specified in the comments.

```
action_space(agent_id) = {
  'Move', # Move 1 tile in any available
    → direction
  'Attack', # Attack an target - with one
    → of 3 styles
  'Target', # Target of agent's action
  'Use', # Use an item
  'Destroy', # Destroy an item
  'Give', # Give an item to a target
  'GiveGold', # Give gold to a target
  'Sell', # Sell an item from market at a
    → given price
  'Buy', # Buy an item from market at a
    → given price
  'Comm' # Communication number of tokens
}
```

Because it is too difficult to deal with multiple arrays contained in multiple dictionaries, we decided to use an extra library "Pufferlib" [8] to vectorize the environment. The resultant observation space is a single array with dimension equal to 26085 for each agent. The vectorized action space is represented by an integer vector with dimension equal to 13. Upper bound of the action vector is defined as [3, 101, 1025, 50, 13, 13, 101, 99, 101, 5, 13, 99, 13]. Interpretation of this upper bound is possible action values plus one

more value to indicate that no action is taken. For example, the 9th dimension indicates the next movement of the agent. Agent can move up, down, to left, and to right. There are 4 possible action values (0, 1, 2, 3). If we enter 4 for this dimension, it means that the agent stays where it is for next game frame.

B. Curriculum Learning

Our reward function is determined by the tasks predefined before training, forming a curriculum. At the beginning of each round, different tasks are randomly generated for different teams. A team receives one point at the end of the round if it accomplishes its assigned task. We assume that the difficulty of tasks is balanced out as teams compete with each other, each with randomly assigned tasks.

It is not easy to organize game environment which involves harvesting food, building armors, trading in market, combatting with NPCs, upgrading professions, and etc. Fortunately, Pufferlib provides a holistic framework for offline agents' emulation, policy pooling, ranking and vectorization. For managing multiple agents' interactions, one intuitive way is to view their roles as 'teachers' and 'students'. An illustration would be the case of asymmetric self-play, where two agents, Alice and Bob, aims for different goals under the same scenario: Alice opens the door with a key then turns on the light, challenges Bob to achieve the same state and Bob attempts to complete it as fast as he can.

Evidently, bad curriculum is worse than no curriculum. Hence, in our case, one has to be attentive when implementing a teacher-guided model for a partially observable Markov decision process to observe unseen states of a 'student'. Recent studies have found Automatic Curriculum Learning performing well for generalist and multi-goal training, hard tasks' solutions [6] and performance improvements on a restricted task set. For baseline, we apply a fixed curriculum of tasks and OpenELM integration. Given \mathcal{H} with any information about past interactions, we learn a task selection function $D : \mathcal{H} \rightarrow \mathcal{T}$.

$$\text{Obj} : \max_D \int_{\mathcal{T} \sim \mathcal{T}_{\text{target}}} P_T^N dT, \quad (1)$$

In the above objective, ACL algorithm maximizes $\frac{P^N}{T}$ over a distribution of target tasks $\mathcal{T}_{\text{target}}$ (i.e., the agent's behavior on task T after N training steps). Instead of learning solutions separately, it would be better to learn a latent skill space so that every task could be represented in a distribution over skills and thus skills are reused among tasks. Therefore, we consider difficulties of different tasks to initialize ACL, which is followed by reward designs aimed at balancing exploration and exploitation. The typical learning process (LP) is discussed in Algorithm 2. Similar to our baseline model of PPO, it alternates between sampling and optimizing. Yet the optimization is not done through stochastic gradient ascent of a surrogate objective function as PPO. Instead, it utilizes the model fitting and selection to guide the sampling process of new parameters for the tasks. As for possible developments, it is probable that an adversarial domain generator can strengthen policies trained for Simulation to Reality applications.

ALGORITHM 2: Learning Process

Require: Student \mathcal{S} , fitting rate N
for timestep= $1, N$ **do**
 Random Sampling $p \in \mathcal{P}$;
 send $E(\tau \sim \mathcal{T}(p))$ to students;
 store its reward and LP;
end
for inner_loops= $1, K$ **do**
 Select Model with best Akaike Information
 Criterion;
 perform epsilon-greedy, send $E(\tau \sim \mathcal{T}(p))$ to
 students, store its reward and LP;
end

IV. RESULTS AND DISCUSSION

In our DQN implementation, we used 13 Q-networks and 13 target Q-networks with 2 hidden layers whose dimension is 256 and 16. We have 128 agents in one environment divided into 16 groups. We performed 100 episodes with maximum timestep equal to 20 for each training. We have 13 replay buffers for 13 action parameters and buffers are shared for all agents. We have two curriculum as mentioned before. A simple one contains only "Staying Alive" task and a comprehensive one contains many different tasks including "Eating Food", "Going Farthest", "Level Up", and so on. All training were performed on a PC Intel i7-10750H, 2.6GHz, 2080s NVIDIA GPU, using a python implementation, which is approximately equivalent to T4 runtime with Google Colab. It takes about 35 minutes for one training.

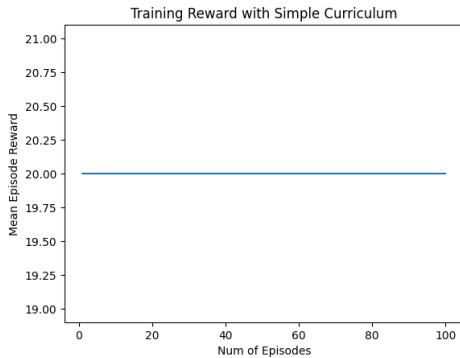


Fig. 1. Mean Episode Reward with the Simple Curriculum

When we use the simple curriculum to train, the mean reward of 128 agents after an episode does not change at all as shown in figure 1. The Q-networks are not updated. This is because the maximum timestep 20 is too small such that no agent dies before 20 game frames. If we insist on training with this simple curriculum, we can imagine that thousands of timesteps are required. To encapsulate more information into the reward function, we changed to the comprehensive curriculum.

With comprehensive curriculum, Q-networks were updated regularly but the mean reward does not converge with a high variance as shown in figure 2. This is the best result we can

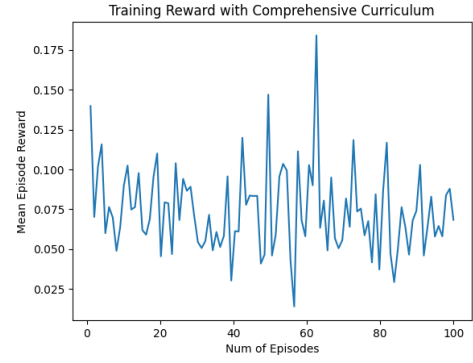


Fig. 2. Mean Episode Reward with the Comprehensive Curriculum

achieve with the resource we have.

But the PPO implementation provided by the Neural MMO group can perform 1K timestep simulation and updates in less than 1 minute. It makes us very interested in why PPO can perform much better than our method. We found following reasons:

1. PPO does not need to swap back and forth between CPU and GPU within episodes. In our implementation, for every game frame, we have to send our states to GPU to evaluate q-values, and send q-values back to CPU to perform greedy policy to determine the action. PPO is a policy-gradient based algorithm. It can perform the simulation with 1000 timesteps first, store the trajectories of all agents, and update the policy once.

2. PPO is more stable because it is updated with episode-based frequency. DQN performs network update within each episode, therefore, its optimization target is moving around within each episode, which is one of the source of high variance. PPO is updated after trajectory is collected after each episode, and it controls the ratio of the update. It is much more stable than DQN when the scale of the environment is huge.

3. PPO makes good use of deep learning strategies. PPO updates its policy using multiple epochs of minibatch updates on each batch of data, whereas DQN performs updates on individual transitions from the replay buffer. Minibatch updates with multiple epochs are well implemented with many libraries.

According to the result, we now understand the advantages of PPO and why it is required for complicated multi-agent environment.

V. CONCLUSIONS

We tried to implement DQN for multi-agent environment, although we did not achieve similar performance compared with PPO. It is recognized that different agents don't act independently, but their decision process and consequences are entwined, and they collaborate or compete for various goals. Whilst this would be challenging to model, we could thereby deploy advanced techniques of curriculum learning to explore progress and exploit synergy, and to finally optimize the survival and gain of agents.

REFERENCES

- [1] Sutton and Barto. Reinforcement Learning, *MIT Press*, 2020.
- [2] Read. Lecture IV - Reinforcement Learning I. In *INF581 Advanced Machine Learning and Autonomous Agents*, 2024.
- [3] Mnih, V., Kavukcuoglu, K., Silver, D. et al. Human-level control through deep reinforcement learning. *Nature* 518, 529–533 (2015). <https://doi.org/10.1038/nature14236>
- [4] John Schulman and Filip Wolski and Prafulla Dhariwal and Alec Radford and Oleg Klimov. Proximal Policy Optimization Algorithms, 2017
- [5] Bick, Daniel (2021) Towards Delivering a Coherent Self-Contained Explanation of Proximal Policy Optimization. Master's Thesis / Essay, Artificial Intelligence.
- [6] Rémy Portelas, Cédric Colas, Lilian Weng, Katja Hofmann, and Pierre-Yves Oudeyer. 2021. Automatic curriculum learning for deep RL: a short survey. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI'20)*. Article 671, 4819–4825.
- [7] Joseph Suarez, Neural MMO 2.0 Documentation link.
- [8] Pufferlib link.