

Learning to play Table Tennis using Multi-agent Reinforcement Learning

Aratrika Basu
University of Southern California
aratrika@usc.edu

Jayantraj Coimbatore Selvakumar
University of Southern California
jcoimbat@usc.edu

Justin Fu
University of Southern California
ziyuefu@usc.edu

Manikanta Chunduru Balaji
University of Southern California
mchundur@usc.edu

Shruthi Ramesh
University of Southern California
shruthir@usc.edu

Sowmya Voona
University of Southern California
voona@usc.edu

Yue Wu
University of Southern California
ywu73061@usc.edu

Abstract—The objective of our project is to understand how we can use machine learning algorithms to train dual AI agents to play a game of table tennis. Through this process we develop both dual-agent and single-agent game, exploring different RL models. Our efforts delivers very capable table-tennis agents that can serve and return similar to a human player. Our models have produced strong performance baselines that should encourage future explorations.

Index Terms—Table Tennis, Game Learning, RL

I. INTRODUCTION

Table Tennis, known as Ping-Pong, is a popular sport in which two or four players hit a lightweight ball on a hard table divided by a net. There has been many attempts to use robots to play table tennis in real-world. But many of those research focus more on perception, such as tracing the trajectories of the ball and predicting its position, transferring the learned agents in simulation environment to real-world setting. They could not compete against human players at any level, either dynamically or strategically. On the other hand, there are also plenty of table tennis games available on virtual environment, like PC, console and mobile devices. However, they are mostly platforms for two human players play against each other, while the provided AI are mostly rule-based. Our table tennis agents aim to achieve high performance under virtual environment, so that learning the strategies are more important to us compared with many prior works on table tennis robots. As a result, we want our environment to be not only neat and clean, but also to consider physical properties like spinning and collision coefficients.

II. RELATED WORKS

A. Reinforcement Learning

Reinforcement learning is a type of machine learning that uses AI systems to follow a policy in order to learn an objective and there by maximize the cumulative reward [1]. Here the AI system starts learning step by step by trial and error approach. For every correct action it performs it is

given a reward and for any subsequent mistake it receives a penalty. Using this feedback mechanism of reward and penalty reinforcement learning learns well in the environment around it [2]. With the development of deep learning, neural networks empower RL with unprecedented abilities in field of Go [?], Atari [4], StarCraft [5], Robotics [6]. A major family of RL algorithms are policy optimization, where they represent a policy as $\pi_{\theta}(a|s)$. The parameters θ are optimized by gradient decent of objectives, often involving learning value functions at the same time. Some representative algorithms are actor-critic algorithm [?], which is temporal difference version of policy gradient, A2C [8], which directly performs gradient ascent in asynchronous manner, and PPO [9], which indirectly maximize a surrogate objective function. Another Family is based no action value function, called Q-Learning, first raised by Watkins in 1992 [10].

B. Table Tennis Robots

Attempts to use robots to play table tennis could be traced back to the 80s. Since Anderson [11] built a real-world vision system which subjectively evaluates and improves its motion plan as the data arrives, many table tennis robot systems were built [12], [13], [14], [15], [16], [17].

As the development of deep learning, especially reinforcement learning, training a robot to play table tennis in real world has been made possible. Lately, Wenbo et al. [18] demonstrate a model-free approach mixed of evolutionary search and CNN-based policy architectures. Jonas et al. [19] shows a modified DDPG [20] could increase sample efficiency in table tennis. Büchler et al., combines step-based reinforcement learning with pneumatic artificial muscles, and achieved great performance using a hybrid sim and real training process. For further learning, Matsushima summarizes the many learning approaches in robotic table tennis [21].

C. multi-Agent learning

Generally, in gaming, it often involves the participation of more than one single agent, which fall into the real of multi-agent RL(MARL). As the previous papers mostly try to solve table-tennis training a single agent, we want to capture the competitive nature of a sport, thus training a pair of agent against each other. MARL algorithms are widely known to be sample-inefficient and millions of interactions are needed. For the game of table tennis, the interaction between the agent and the environment is relatively simple compared with games like starcraft. Hence, we would focus more on the model-free setting, where the policies are learned without direct access to the environment.

Compared with single-agent RL, MARL suffers from several challenges. As summarized by [22], MARL does not have unique learning goals and whether convergence of equilibrium point is the alpha performance criterion for MARL algorithm analysis is controversial. Some researchers found value-based MARL algorithms fail to converge to stationary Nash equilibrium point for general-sum Markov games [23]. Another major issue is the non-stationary setting as multiple agents could simultaneously interact with the environment and each other. This could bring challenge to value estimation as well as policy optimization during training. Scalability is a issue coming along with non-stationary, as the joint action space is exponentially increasing. Even in a dual agent setting as table tennis, the sample efficiency would still be a major bottleneck.

MARL has many information structures(who knows what at the training and execution) [22]. For the dual agent setting of table tennis, the straight forward way is treat other agent as part of the environment, which is called *Independent Learning(IL)*. But IL face the problem of non-stationary dynamic, which harms the performance of policies. Some work try to stabilize the learning process [24], [25]. Others try to build communication protocols between agents [26], [27]. Another major MARL learning diagram is *Centralized Training and Decentralized Execution (CTDE)*. One Representative CTDE method is MADDPG [28], a multi-agent version actor-critic. Each agent maintains its own critic Q_i , which estimates the joint value function and uses the critic to update its decentralized policy.

MARL consists of three groups, fully cooperative, fully competitive and mixed of two. Though Table tennis is considered to be competitive game, we would also try some mixed methods since table-tennis is not a typical *zero-sum* game, where the reward for one player is exactly the loss of the other.

III. DATA AND ENVIRONMENTS

A major part of this project is building an environment that is interactive, visual attractive, and efficient for training artificial agents.

A. Development Tools

To develop the environment, we use the following tools.

- 1) **Unity 3D: 2020.3.20** Unity is a cross-platform game engine developed by Unity Technologies. The game engine can be used to develop interactive 3D, 2D, as well as interactive simulations and other experiences [?]. Unity version 2020.3.20 is utilized for the environment setup.
- 2) **Unity Machine Learning Agents Toolkit** The Unity Machine Learning Agents Toolkit (ML-Agents) [29] is an open-source project that enables games and simulations to serve as environments for training intelligent agents. They provide state-of-the-art algorithms which can be used to train intelligent agents to play different 3D and 2D games. The ML agents package provides an option to convert a Unity scene into a learning environment where character behaviors can be trained using machine learning algorithms.
- 3) **Pytorch** PyTorch is an open-source machine learning library based on the Torch library, used for applications such as computer vision and natural language processing.
- 4) **Python** The most popular language for machine learning research.
- 5) **Tensorboard** It is a Tensorflow visualization toolkit that provides visualization and tools for machine learning experimentation. It helps to track and visualizing metrics like loss and accuracy. Tensorflow.dev provides an easy way to share ML experimentation results.

B. Environment Functionality

Our project aimed to create an environment and ML agents, to enable them to play a game of table tennis. For this purpose, we found out that Unity provides a comprehensive environment where game objects can be created and modeled as per user requirements. Also, game objects can be used as ML agents and can be trained using Proximal Policy Optimization and Soft Actor-Critic model provided by Unity. Therefore, we selected the Unity platform as the environment.

Our environment contains a Table Tennis bat with the ability to assign different unity materials to different sections of the bat for customizability. The sections include Bat forehand face, Bat backhand face, Bat center, and the Bat handle. The second model in this pack is the table tennis table which can have different unity materials assigned to it for customizability. The sections include Tabletop, Table Legs, Net frame.

- 1) **Vector Observations:** From the environment, we are collecting the positions of bat A, bat B, and the ball. Also, we are collecting the velocity of bat A, bat B, and the ball. We have used these observations to train the model using different Reinforcement learning algorithms.
- 2) **Actions:** We have designed the environment in a way where the bats can move along X and Y axes and can rotate along X axes. Our goal is also to use the Z axes in the following weeks. The bats can also be moved using the 'right', 'left' keys.

3) **Reward Policy:** We have designed our reward policy in such a way that if a player commits a mistake or makes a foul move the opponent player gets the reward for it. The following are the foul moves implemented for our project:

- Player hitting the ball to the net.
- Player hitting the ball over the boundary.
- Ball bouncing more than once on the same side of the court.

For implementing the reward policy we are keeping track of the parameters given below:

- Last Hit Agent : The agent who hit the ball previously before coming to the current player.
- Last Collided With : This keeps track of the last surface the ball collided with. Here the surface refers to the court of player A, court of player B, Net, etc.
- Next Agent Turn: This keeps the track of the next agent who has to hit the ball to continue the game.

Using the above parameters we are rewarding the agents.

The following figure shows the environment we have built for now.

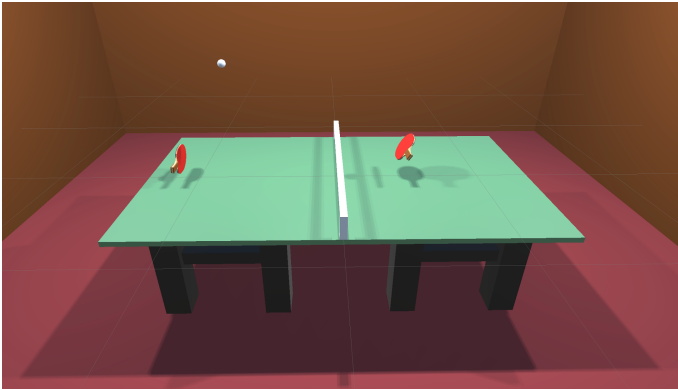


Fig. 1: Table Tennis Environment

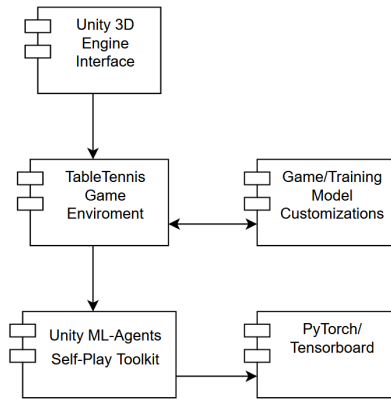


Fig. 2: System Overview

C. Software Overview

The final software we build consists of the following parts.

1) **Game Controller Class:** This is the main controller class that is interlinked to all other classes. It has the following functionalities:

- agentScores(): This method is used for rewarding the agents.
- episodeReset(): It is used to reset the episode for every foul move.
- matchReset(): It is used to reset the match.
- ballHitsAgent(), ballHitsFloor(), ballHitsBoundary(), agentHitsNet(), ballHitReward(): These methods are used to handle the reward for the agent depending on the foul moves described for the game.

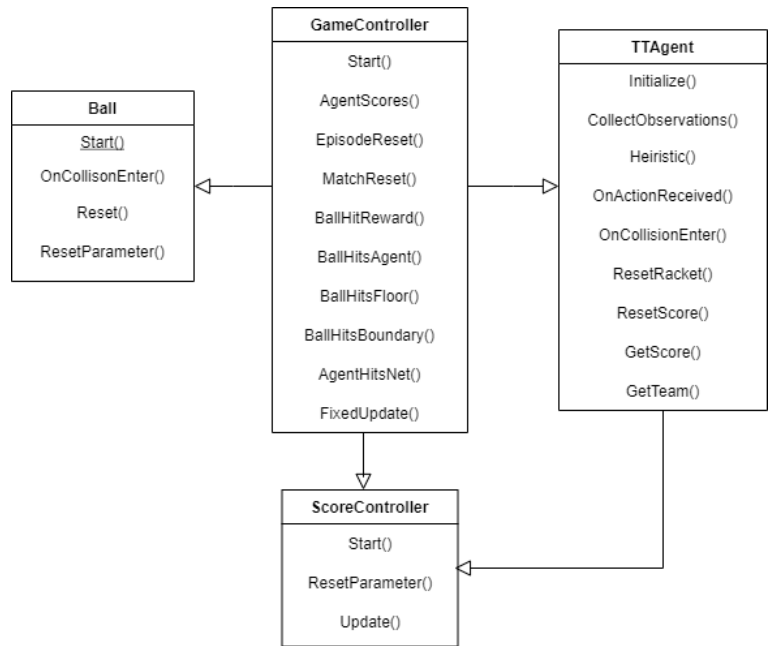


Fig. 3: Software Workflow Diagram

2) **Ball Class:** This class refers to the functionalities used for the ball.

- onCollisionEnter(): This function handles the different nuances when the ball collides with different surfaces. For example, when the ball collides with bat A we are checking the parameter status of the lastHitAgent, lastCollidedWith and we are rewarding the agent as per the rules of table tennis.
- reset(), resetParameter(): This functionality handles resetting the ball positions for every episode.

3) **TAgent Class:** This class involves all the functionalities required for the agent.

- CollectObservations(): It is used for collecting the vector observations for the bats and the ball. The velocity observations of the bats and the ball are also collected.

- `Heuristic()`: This functionality is used to assign the movement to the bat along 'x' and 'y' axis. The bat can be moved along the horizontal axis using the right and left keys. It can be moved in the vertical direction using the upward and downward keys. The racket can be made to jump using the 'X' key.
- `OnActionReceived`: This executes the actions by moving the game objects in the vector space.
- `resetRacket`: It is used to reset the racket position.
- `resetScore`: It is used to reset the score for the agents.

4) Score Controller Class

IV. METHODS

A. Preliminaries

The Table Tennis Game could be described as a Markov Process, and is a Markov Game [30].

Markov Decision Process(MDP). An MDP is defined as

$$\langle S, A, T, R, \rho, \lambda \rangle$$

where S a set of states, A a set of actions, $T : S \times A \rightarrow P(S)$ a stochastic transition function, $R : S \times A \rightarrow R$ a reward function, $\lambda \in [0, 1)$ a discount factor. The agent(table tennis player) interacts with the ball by performing its policy $\pi : S \rightarrow P(A)$. The agents learn this policy to maximize the expected cumulative discounted reward:

$$J(\pi) = E_{\rho, \pi, T} \sum_{t=0}^{\infty} r_t \lambda^t$$

where $r_t = R(s_t, a_t)$, $s_0 \sim \rho_0(s_0)$, $a_t \sim \pi(s_t)$, $s_{t+1} \sim T(\cdot | s_t, a_t)$

Markov Game(MG). An MG is an extension of MDP and is defined as

$$\langle S, N, \{A^i\}_{i=1}^N, \{R^i\}_{i=1}^N, \{O^i\}_{i=1}^N, \rho, \lambda, Z \rangle$$

where the action sets now contain N agents, namely, $A^1 \dots A^N$, state transition function $T : S \times A^1 \dots A^N \rightarrow P(S)$, reward function $R : S \times A^1 \dots A^N \rightarrow R$. For partially observable Markov games, each agent i receives local observation $o^i : Z(S, i) \rightarrow O^i$ and interacts with environment with its policy $\pi^i : O^i \rightarrow P(A^i)$. The expected cumulative discount reward now turns into

$$J^i(\pi^i) = E_{\rho, \pi^1, \dots, \pi^N, T} \sum_{t=0}^{\infty} r_t^i \lambda^t$$

where $r_t^i = R^i(s_t, a_t^1, \dots, a_t^N)$. we will discuss methods like PPO, SAC and DQN and how we train our table tennis agent using these RL algorithms under Markov Game setting.

B. Models

1) *Proximal Policy Optimization(PPO)*: Proximal Policy Optimization(PPO) is an on-policy based reinforcement learning algorithm. This algorithm was introduced by the OpenAI team in the year 2017 [9] and quickly became one of the most popular RL methods surpassing the Deep-Q learning method.

PPO is scalable, data efficient, and successful on a variety of problems without hyper-parameter tuning.

PPO is an algorithm that attains the data efficiency and reliable performance of trust region policy optimization (TRPO), while using only first-order optimization. It involves collecting a small batch of experiences interacting with the environment and using that batch to update its decision-making policy. Once the policy is updated with this batch, the experiences are thrown away and a newer batch is collected with the newly updated policy. This is the reason it is an ‘‘on-policy learning’’ approach where the experience samples collected are only useful for updating the current policy once.

PPO improves stability of the learning by mainly 2 techniques:

- **Clipped Surrogate Objective**: The Clipped Surrogate Objective is a drop-in replacement for the policy gradient objective that is designed to improve training stability by limiting the change you make to your policy at each step.
- **Multiple epochs for policy updating** : Unlike vanilla policy gradient methods, and because of the Clipped Surrogate Objective function, PPO allows user to run multiple epochs of gradient ascent on your samples without causing destructively large policy updates. This allows to squeeze more out of your data and reduce sample inefficiency.

Algorithm 1 PPO, Actor-Critic Style

```

for iteration=1,2,... do
  for actor=1,2,...,N do
    Run policy  $\pi_{\theta, \text{old}}$  in environment for  $T$  timesteps
    Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
  end for
  Optimize surrogate  $L$  wrt  $\theta$ , with  $K$  epochs and minibatch size  $M \leq NT$ 
   $\theta_{\text{old}} \leftarrow \theta$ 
end for

```

Fig. 4: Proximal Policy Optimization Algorithm [9]

2) *Soft Actor Critic(SAC)*: Soft Actor Critic(SAC) is an off-policy model-free reinforcement learning algorithm. This RL algorithm was developed jointly by UC Berkely and Google and was introduced in the year 2018 [31]. It is considered one of the most efficient algorithm to be used in real-world robotics.

The biggest feature of SAC is that it uses a modified RL objective function. Instead of only seeking to maximize the lifetime rewards, SAC seeks to also maximize the entropy of the policy. A high entropy in our policy explicitly encourages exploration, encourages the policy to assign equal probabilities to actions that have same or nearly equal Q-values, and also ensures that it does not collapse into repeatedly selecting a particular action that could exploit some inconsistency in the approximated Q function. SAC overcomes the brittleness problem by encouraging the policy network to explore and not assign a very high probability to any one part of the range of actions.

3) *DQN*: DQN is an off-policy, value-based, model-free RL algorithm. This algorithm was introduced by DeepMind Technologies in the year 2013 [32]. The algorithm was modified in the 2015.

Algorithm 1 Soft Actor-Critic

Initialize parameter vectors $\psi, \bar{\psi}, \theta, \phi$.
for each iteration **do**
 for each environment step **do**
 $\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t|\mathbf{s}_t)$
 $\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$
 $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$
 end for
 for each gradient step **do**
 $\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$
 $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$
 $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$
 $\bar{\psi} \leftarrow \tau \bar{\psi} + (1 - \tau)\psi$
 end for
end for

Fig. 5: Soft Actor-Critic Algorithm [25]

A Deep Q-Network, based on Q-learning framework, uses a deep network to approximate the state-value pair. In some games, it takes in several frames of the game as input and returns state values for each action as outputs. DQN comes up with Experience Replay, an idea to store the episode steps in memory of off-policy learning, where samples are drawn from the replay memory at random. Having a replay memory makes the problem more like a supervised learning problem. Another difference to normal Q-learning is that Q-Network periodically updated with the latest weights and optimized towards a frozen target network in every steps. The former helps training to be more stable since it prevents short-term oscillations from a moving target. The later deals with autocorrelation that would occur from on-line learning,

DQN overcomes unstable learning by mainly 2 techniques.

- Experience Replay
- Target Network

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N
Initialize action-value function Q with random weights
for episode = 1, M **do**
 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
 for $t = 1, T$ **do**
 With probability ϵ select a random action a_t
 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
 Execute action a_t in emulator and observe reward r_t and image x_{t+1}
 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}
 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
 end for
end for

Fig. 6: DQN Algorithm [32]

4) *MA-POCA*: MultiAgent POsthumous Credit Assignment [29] is a novel multi-agent trainer that trains a centralized critic, a neural network that acts as a "coach" for a whole group of agents. Rewards can be given to the team as a whole, and the agents will learn the best ways to contribute to achieving that reward. Agents can also be given rewards individually, and the team will work together to help the individual achieve those goals.

Additionally in MA-POCA agents can be added or removed from the group during an episode, such as when agents spawn or die in a game. If agents are removed mid-episode, they will still learn whether their actions contributed to the team winning later. This enables the agents to take group-beneficial actions even if it results in them being removed from the game. MA-POCA can also be combined with self-play to train teams of agents to play against each other

5) *DDPG*: Deep Deterministic Policy Gradient (DDPG) is an off-policy algorithm which is a version of Q-learning for continuous action space. It combines Q-learning and gradient descent in that it uses the Bellman equation to learn the Q-function, and uses the Q-function to learn the policy

Given optimal action-value function $Q^*(s, a)$, then in any given state, the optimal action $a^*(s)$ can be found by solving

$$a^*(s) = \arg \max_a Q^*(s, a)$$

DDPG interleaves learning an approximator to $Q^*(s, a)$ with learning an approximator to $a^*(s)$

6) *Curriculum Learning*: Curriculum learning is a learning agenda to progressively learn from simple to hard circumstances. The idea to imitate human's learning progress under curriculum could be traced back to as early as 1993, when Jeffery Elman proposed a strategy to begin training neural networks with a restricted set of simple data and graduate expand to complex training samples.

For our table tennis game we make the following curriculum for the reward

- 1) Training the agent with initially high reward values set for hitting the ball regardless of foul moves and low reward values for scoring against opponent.
- 2) As the agent progresses through lessons, reward for hitting the ball is reduced and reward for scoring against opponent is increased

V. RESULTS AND ANALYSIS

A. PPO

The PPO training is carried out by setting the behavior type of agents to "Default" in Unity so that no external/human interaction is required to play the game. We used the mlagents-learn package to execute the configuration file which contains the hyper parameters specific to each model. Each time a configuration file is called a new model is trained and gets saved in the local system. Later, the trained model can be embedded into Unity as the model type in order to observe the learning that the agents have obtained.

We have tuned the model by using a variety of hyper parameter combination in our configuration file while keeping our batch size as 2048, hidden units in each layer as 256 and initial ELO as 1200. The table [I] contains the hyper parameter combinations.

Buffer Size	Lr	Epochs	Schedule	Layers	Steps	Final ELO
2048000	0.0003	3	constant	3	370000	1202
2048000	0.0003	3	constant	2	1.6M	1191
2048000	0.001	3	constant	2	1.93M	1208
20480	0.03	3	constant	2	730000	1203
20480	0.01	3	constant	2	2.23M	1190
20480	0.01	3	constant	3	2.19M	1170
20480	0.01	500	linear	3	20000	1193
20480	0.01	10	linear	3	1.2M	1189
20480	0.01	1000	linear	3	100000	1195
20480	0.01	100	linear	3	1M	1205

TABLE I: PPO Hyper parameter Combination

B. SAC

The training is carried out by setting the behavior type of agents to "Default" in Unity so that no external/human interaction is required to play the game. We used the mlagents-learn package to execute the configuration file which contains the hyper parameters specific to each model. Each time a configuration file is called a new model is trained and gets saved in the local system. Later, the trained model can be embedded into Unity as the model type in order to observe the learning that the agents have obtained.

We have tuned the model by using a variety of hyper parameter combination in our configuration file while keeping our hidden units in each layer as 256, learning rate schedule as 'constant' and initial ELO as 1200. The table [II] contains the hyper parameter combinations.

Buffer size	Batch size	Learn Rate	init steps	Bounce	Layer	Steps	Final ELO
500000	128	0.0003	0	1	2	8M	2352
50000	128	0.01	0	1	2	2M	1272
500000	128	0.003	0	1	2	3.6M	1540
1000000	1024	0.0003	1000	1	2	4M	2002
500000	512	0.0003	1000	1	2	3M	1953
500000	512	0.0003	1000	1	2	10M	2130
1000000	1024	0.0003	1000	1	3	3.9M	1915
1000000	1024	0.0003	1000	1	3	0.7M	1748
500000	512	0.003	0	0.9	3	1.98M	2005
500000	512	0.0003	0	0.9	2	6.8M	1940

TABLE II: SAC Hyper parameter Combination

C. MA-POCA

We have tuned the model by using a variety of hyper parameter combination in our configuration file while keeping our hidden units in each layer as 256, learning rate schedule as 'constant' and initial ELO as 1200. The table [III] contains the hyper parameter combinations.

Buffer size	Batch	Lr	Hidden	Layers	Steps	Final ELO
20480	2048	0.0003	512	2	1000(A)/4000(B)	1212
20480	2048	0.003	512	2	1000(A)/4000(B)	1270

TABLE III: MA-POCA Hyper parameter Combination

D. Curriculum Learning

Training:

- 1) We train the agent with initial high reward values assigned for hitting the ball, regardless of whether the

move made was a foul, and initial low reward values assigned for scoring against the opponent. But as the agent continues to move through lessons, we progressively reduce the reward value assigned for hitting the ball and increase the reward value set for scoring against the opponent in order to teach the agent to play legal moves with a higher scoring probability. We used progress, represented by the ratio of current steps to maximum steps, as the measure to shift through lessons in our curricula.

- 2) We train the agent by giving it an easy level of initial play, with limited movement required to hit the ball, by setting the initial bat size to a large value. We progressively increase the difficulty through the sequence of lessons, by reducing the bat size (in the z and y axis in Unity) until it reaches the pre-established normal bat size value in the final lesson. This is done in order to teach the agent to increase its range of motion in order to hit the ball. We used progress, represented by the ratio of current steps to maximum steps, as the measure to shift through lessons in our curricula.

E. Single Agent Setting

From the sections before, we have seen that the applying RL methods in a competitive dual agent setting could help agents to learn play the table tennis game. However, we found that some of the methods like PPO could not guide the agents to play the game well. The dual agent setup also restricted us to training the agents with the models provided by unity as we wanted to explore how the agents perform when trained with other different models from external frameworks like gym-library.

We are also curious if we could train a single agent, who only serves the purpose of hitting the ball on to the other side of the court abiding Table Tennis Rules, could learn to play the table tennis game. So we consider the following scenarios:

Environment

- 1) **Agent:** The agent has the same functionalities built-in as in the dual agent setup, but the reward policies are now only set based on how well the agent is able to hit the ball and abide by the table tennis rules.
- 2) **Serve Bot:** A serving bot that serves the ball to the agent with different levels of difficulty depending on the mode it is set to.
- 3) **Reward Policy:**

a) Positive Rewards:

- i) A positive reward is added to the agent whenever it is observed that the ball has collided with the agent through a trigger event
- ii) A positive reward is added to the agent whenever it is able to hit the ball across the net, by placing an invisible object on top of the net to allow us to observe the event whenever the ball passes through this object.

- iii) A positive reward when the agent successfully hits the ball not only across the net but also onto the opponent's table

b) **Penalties:**

- i) when it misses the ball
- ii) when hits the ball twice on its turn
- iii) when hits the ball twice on its own side of the table
- iv) when hits the ball directly onto any boundary or net.

Initial Conditions

- a) **Basic Setup:** Fix the velocity and height from which the ball is served. The ball position on the X and Z axis is randomized. Hence the only challenge that the agent faces is to move towards the ball and be able to hit it.
- b) **Randomized velocity:** In this mode, the agent is challenged with the ball being served with different values for velocities from randomized X, Y, Z positions. The agent has to learn to respond accordingly.
- c) **Randomized Spin:** In this setup, we have added angular velocity and torque to the ball as it is served from the serve bot. The angular velocity and the torque is randomized for every time the serve bot serves the ball. The inclusion of angular velocity and torque results in the ball moving in different directions after bouncing on the table.

The training is carried out by setting the behavior of the single agent to "Default" in Unity so that no external action is required to play the game. The serve bot is here considered as the other agent for our training purpose. The trained model, including SAC, PPO and curriculum learning, can be embedded into Unity single agent to observe the performance of the agent.

F. Analysis

In this section, we would discuss what we have observed and thought on some of our training results.

1) *RL models: SAC, PPO, M-POCA* : Next Figure shows the training result with different methods like PPO, SAC, MA-POCA, etc. In self play as we know ELO is the most important factor that dominates the performance of each model. In our case the SAC model reaches the maximum ELO of 2352 after training for 8M steps and continues to grow after. We have used a very low learning rate of 0.0003 so that the model can learn slowly but efficiently over a longer period of time. In our experiments we have trained multiple SAC and PPO models while varying the parameters and achieved the highest performance on the SAC model having a batch size of 128, learning rate of 0.0003, buffer initial steps of 0, ball bounce of 1 and 2 neural net layers, running for 8M steps.

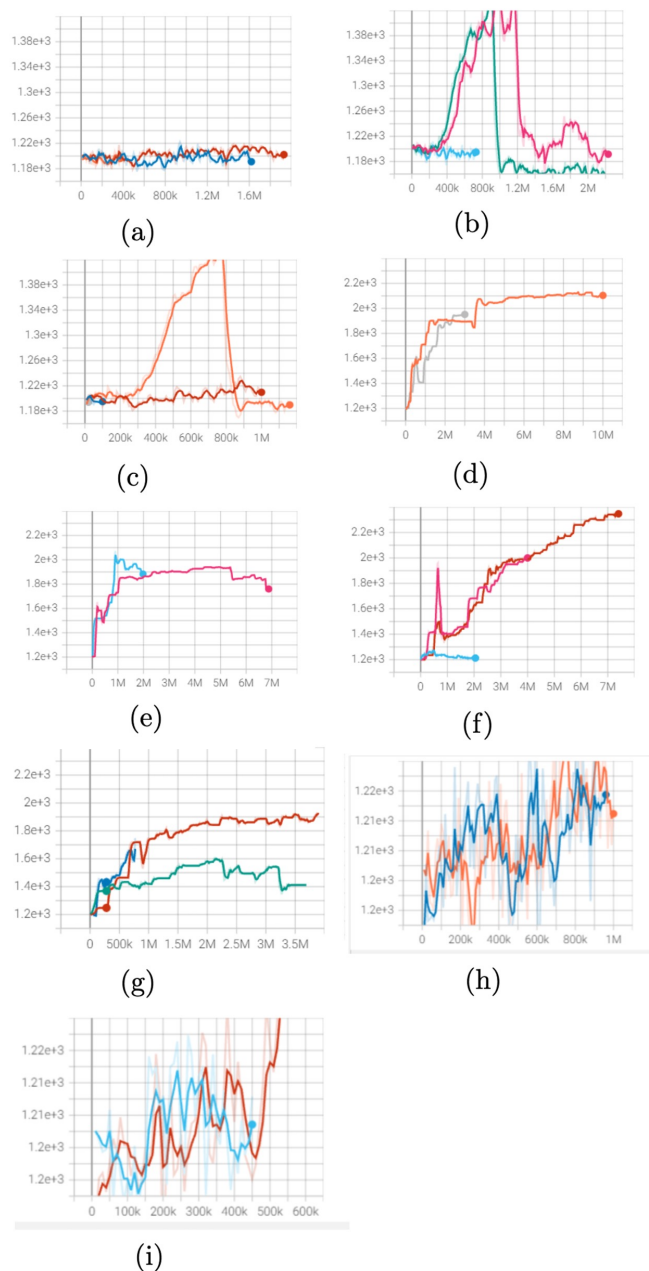


Fig. 7: **PPO training Observation:** Self Play/ELO graphs with the hyper parameter combination:(a) High Buffer Size of 2048000, low Learning Rate of 0.0003, (b) Low Buffer Size of 20480, high Learning Rate of 0.01, epochs of 3, (c) Low Buffer Size of 20480, high Learning Rate of 0.01, epochs of 500, 10, 1000, 100. **SAC training Observation:** Self Play/ELO graphs with the hyper parameter combination:(d) Medium Batch Size of 512, low Learning Rate of 0.0003, buffer initial steps of 1000,ball bounce of 1 with 2 layers (e) Medium Batch Size of 512, low Learning Rate of 0.0003, buffer initial steps of 0,ball bounce of 0.9 with 3 layers (f) Low Batch Size of 128, high Learning Rate varied between 0.01, 0.0003, 0.003, buffer initial steps of 0 with 2 layers (g) High batch Size of 20480, low Learning Rate of 0.0003,buffer initial steps of 1000 experimented with 2 and 3 layers **POCA training Observation:** Self Play/ELO graphs with the hyper parameter combination: (h)High Batch Size of 2048, low Learning Rate of 0.0003, swap steps of 1000 for Agent A and 4000 for Agent B with 2 neural layers (i)High Batch Size of 2048, high Learning Rate of 0.003, swap steps of 1000 for Agent A and 4000 for Agent B with 2 neural layers

REFERENCES

- [1] Wikipedia contributors, “Reinforcement learning — Wikipedia, the free encyclopedia,” https://en.wikipedia.org/w/index.php?title=Reinforcement_learning&oldid=1042314112, 2021, [Online; accessed 6-September-2021].
- [2] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. The MIT Press, 2018. [Online]. Available: <http://incompleteideas.net/book/the-book-2nd.html>
- [3] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, “Mastering the game of Go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, Jan. 2016.
- [4] M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. Azar, and D. Silver, “Rainbow: Combining improvements in deep reinforcement learning,” 2017, cite arxiv:1710.02298Comment: Under review as a conference paper at AAAI 2018. [Online]. Available: <http://arxiv.org/abs/1710.02298>
- [5] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, J. Oh, D. Horgan, M. Kroiss, I. Danihelka, A. Huang, L. Sifre, T. Cai, J. P. Agapiou, M. Jaderberg, A. S. Vezhnevets, R. Leblond, T. Pohlen, V. Dalibard, D. Budden, Y. Sulsky, J. Molloy, T. L. Paine, Gülgehr, Z. Wang, T. Pfaff, Y. Wu, R. Ring, D. Yogatama, D. Wünsch, K. McKinney, O. Smith, T. Schaul, T. P. Lillicrap, K. Kavukcuoglu, D. Hassabis, C. Apps, and D. Silver, “Grandmaster level in starcraft ii using multi-agent reinforcement learning,” *Nat.*, vol. 575, no. 7782, pp. 350–354, 2019. [Online]. Available: <http://dblp.uni-trier.de/db/journals/nature/nature575.html#VinyalsBCMDCCPE19>
- [6] OpenAI, M. Andrychowicz, B. Baker, M. Chociej, R. Józefowicz, B. McGrew, J. W. Pachocki, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba, “Learning dexterous in-hand manipulation,” *CoRR*, vol. abs/1808.00177, 2018. [Online]. Available: <http://dblp.uni-trier.de/db/journals/corr/corr1808.html#abs-1808-00177>
- [7] V. Konda and J. Tsitsiklis, “Actor-critic algorithms,” in *Advances in Neural Information Processing Systems*, S.olla, T. Leen, and K. Müller, Eds., vol. 12. MIT Press, 2000. [Online]. Available: <https://proceedings.neurips.cc/paper/1999/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf>
- [8] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” in *Proceedings of The 33rd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. F. Balcan and K. Q. Weinberger, Eds., vol. 48. New York, New York, USA: PMLR, 20–22 Jun 2016, pp. 1928–1937. [Online]. Available: <https://proceedings.mlr.press/v48/mniha16.html>
- [9] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *CoRR*, vol. abs/1707.06347, 2017. [Online]. Available: <http://arxiv.org/abs/1707.06347>
- [10] C. J. C. H. Watkins and P. Dayan, “Technical note q-learning,” *Mach. Learn.*, vol. 8, pp. 279–292, 1992. [Online]. Available: <https://doi.org/10.1007/BF00992698>
- [11] R. L. Anderson, *A Robot Ping-Pong Player: Experiment in Real-Time Intelligent Control*. Cambridge, MA, USA: MIT Press, 1988.
- [12] F. Miyazaki, M. Matsushima, and M. Takeuchi, “Learning to dynamically manipulate: A table tennis robot controls a ball and rallies with a human being,” *Advances in Robot Control: From Everyday Physics to Human-Like Movements*, 01 2006.
- [13] K. Mülling and J. Peters, *A Computational Model of Human Table Tennis for Robot Application*, 2009, p. 57.
- [14] K. Muelling, J. Kober, and J. Peters, “Learning table tennis with a mixture of motor primitives,” in *2010 10th IEEE-RAS International Conference on Humanoid Robots*, 2010, pp. 411–416.
- [15] Y. Huang, D. Buchler, O. Koc, B. Schölkopf, and J. Peters, “Jointly learning trajectory generation and hitting point prediction in robot table tennis,” in *16th IEEE-RAS International Conference on Humanoid Robots, Humanoids 2016, Cancun, Mexico, November 15-17, 2016*. IEEE, 2016, pp. 650–655. [Online]. Available: <https://doi.org/10.1109/HUMANOIDS.2016.7803343>
- [16] R. Mahjourian, N. Jaitly, N. Lazic, S. Levine, and R. Miikkulainen, “Hierarchical policy design for sample-efficient learning of robot table tennis through self-play,” *CoRR*, vol. abs/1811.12927, 2018. [Online]. Available: <http://dblp.uni-trier.de/db/journals/corr/corr1811.html#abs-1811-12927>
- [17] K. Muelling, J. Kober, O. Kroemer, and J. Peters, “Learning to select and generalize striking movements in robot table tennis,” in *AAAI Fall Symposium on Robots that Learn Interactively from Human Teachers*, 2012, pp. 263–279. [Online]. Available: <http://www.aaai.org/ocs/index.php/FSS/FSS12/paper/view/5602>
- [18] W. Gao, L. Graesser, K. Choromanski, X. Song, N. Lazic, P. Sanketi, V. Sindhwani, and N. Jaitly, “Robotic table tennis with model-free reinforcement learning,” 2020.
- [19] J. Tebbe, L. Krauch, Y. Gao, and A. Zell, “Sample-efficient reinforcement learning in robotic table tennis,” 2021.
- [20] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” in *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2016. [Online]. Available: <http://arxiv.org/abs/1509.02971>
- [21] M. Matsushima, T. Hashimoto, M. Takeuchi, and F. Miyazaki, “A learning approach to robotic table tennis,” *IEEE Transactions on Robotics*, vol. 21, no. 4, pp. 767–771, 2005.
- [22] K. Zhang, Z. Yang, and T. Başar, “Multi-agent reinforcement learning: A selective overview of theories and algorithms,” 2021.
- [23] Y. Shoham, R. Powers, and T. Grenager, “Multi-agent reinforcement learning: a critical survey,” *Tech. Rep.*, 2003.
- [24] G. Palmer, K. Tuyls, D. Bloembergen, and R. Savani, “Lenient multi-agent deep reinforcement learning,” *CoRR*, vol. abs/1707.04402, 2017. [Online]. Available: <http://dblp.uni-trier.de/db/journals/corr/corr1707.html#PalmerTBS17>
- [25] S. Sukhbaatar, A. Szlam, and R. Fergus, “Learning multiagent communication with backpropagation,” 2016.
- [26] X. Kong, B. Xin, F. Liu, and Y. Wang, “Revisiting the master-slave architecture in multi-agent deep reinforcement learning,” 2017.
- [27] J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson, “Learning to communicate with deep multi-agent reinforcement learning,” in *NIPS*, D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, Eds., 2016, pp. 2137–2145. [Online]. Available: <http://dblp.uni-trier.de/db/conf/nips/nips2016.html#FoersterAFW16>
- [28] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, “Multi-agent actor-critic for mixed cooperative-competitive environments,” in *NIPS*, I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, Eds., 2017, pp. 6379–6390. [Online]. Available: <http://dblp.uni-trier.de/db/conf/nips/nips2017.html#LoweWTHAM17>
- [29] A. Juliani, V.-P. Berges, E. Teng, A. Cohen, J. Harper, C. Elion, C. Goy, Y. Gao, H. Henry, M. Mattar, and D. Lange, “Unity: A general platform for intelligent agents,” 2020.
- [30] M. L. Littman, “Value-function reinforcement learning in markov games,” *Cogn. Syst. Res.*, vol. 2, no. 1, pp. 55–66, 2001. [Online]. Available: <http://dblp.uni-trier.de/db/journals/cogsr/cogsr2.html#Littman01>
- [31] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, ser. Proceedings of Machine Learning Research, J. G. Dy and A. Krause, Eds., vol. 80. PMLR, 2018, pp. 1856–1865. [Online]. Available: <http://proceedings.mlr.press/v80/haarnoja18b.html>
- [32] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. A. Riedmiller, “Playing atari with deep reinforcement learning,” *CoRR*, vol. abs/1312.5602, 2013. [Online]. Available: <http://arxiv.org/abs/1312.5602>