

Optimal Stroke Learning with Policy Gradient Approach for Robotic Table Tennis

Yapeng Gao^{1*}, Jonas Tebbe¹ and Andreas Zell¹

^{1*}Cognitive Systems, Eberhard Karls University Tübingen, Geschwister-Scholl-Platz, Tübingen, 72074, Baden-Württemberg, Germany.

*Corresponding author(s). E-mail(s):

yapeng.gao@uni-tuebingen.de;

Contributing authors: jonas.tebbe@uni-tuebingen.de;

andreas.zell@uni-tuebingen.de;

Abstract

Learning to play table tennis is a challenging task for robots, as a wide variety of strokes required. Recent advances have shown that deep Reinforcement Learning (RL) is able to successfully learn the optimal actions in a simulated environment. However, the applicability of RL in real scenarios remains limited due to the high exploration effort. In this work, we propose a realistic simulation environment in which multiple models are built for the dynamics of the ball and the kinematics of the robot. Instead of training an end-to-end RL model, a novel policy gradient approach with TD3 backbone is proposed to learn the racket strokes based on the predicted state of the ball at the hitting time. In the experiments, we show that the proposed approach significantly outperforms the existing RL methods in simulation. Furthermore, to cross the domain from simulation to reality, we adopt an efficient retraining method and test it in three real scenarios. The resulting success rate is 98% and the distance error is around 24.9 cm. The total training time is about 1.5 hours.

Keywords: Table tennis robot, Stroke learning, Reinforcement learning, Sim2Real

1 Introduction

Reinforcement learning (RL) [1, 2] is a general learning paradigm that addresses the problem of how an acting agent can learn an optimal behavioral strategy while interacting with an unknown environment surrounding it [3]. Recently, RL has achieved a variety of successes, most prominently in autonomous driving [4, 5], gaming [6, 7], and robotics manipulation [8, 9]. Notably, most RL research is based on the formalism of Markov decision process (MDP) which models the decision-making problem. Specifically, while interacting with the environment, the agent takes actions based on the observed state of the environment and then receives rewards for its actions. Thus, the goal of RL is to maximize the expected value of the cumulative reward in an episode. However, formulating a near-optimal policy for the agent's actions usually requires extensive exploration in the action space. For example, OpenAI Five [7] defeated the world champions at an esports (Dota2) game by developing a training system using RL techniques. To fully explore the action space, it trains for 180 days on 256 GPUs and 128,000 CPU cores, based on an average of 180 days worth of self-play. Therefore, training RL models in the context of robotics is particularly challenging, since in such a context it is very difficult to safely collect samples that cover all possible actions in their space.

A common method to address this problem is training RL models with the help of simulations [10–14]. Previous work has shown that simulation can be used as a valuable tool for robotics research, as performing robotic skills in simulation is comparatively easier than in the real world [5, 15]. The use of simulation greatly facilitates the implementation of RL in robotics by allowing a comprehensive exploration of the robotic action space with less engineering effort than in real world by adjusting the parameters of the simulator. However, policies learned in simulation are often unsuitable for reality due to the reality gap. To bridge the gap, researchers usually incorporate data from the real world for training, or retrain the models in reality. For example, a table tennis robot was trained with RL in a hybrid simulation and real system [16]. The real trajectories of the ball were recorded and replayed in simulation to use the real data as much as possible.

However, one problem in existing RL methods such as Trust Region Policy Optimization (TRPO) [17], Proximal Policy Optimization (PPO) [18], Deep Deterministic Policy Gradient (DDPG) [19], Twin Delayed DDPG (TD3) [20], or Soft Actor-Critic (SAC) [21] is that the fuzzy one-dimensional reward cannot precisely express the interaction with the environment for multi-dimensional actions. Therefore, inspired by the previous work, we propose a novel approach for optimal stroke learning in robotic table tennis. A 3D Q -value function is designed to cope with the corresponding 3D reward vector. Two learning steps, including training in simulation and retraining in reality, can be completed in around 1.5 hours for the balls with a wide variety of spins, speeds, and positions. The main contributions of this work are as follows:

- We design a realistic table tennis robot simulation system for optimal stroke learning with RL, as shown in Fig. 1 left. The simulation is based on the Gazebo simulator, the Robot Operating System (ROS), OpenAI Gym [22], and the RL library Spinning Up. Controlled by the Gazebo plugin, the robot and the table tennis ball can publish their states via ROS topics.
- We decompose the learning strategy into two stages: first, the prediction of the ball's hitting state, and second, the learning of the optimal stroke, which is the focus of this paper. Based on the controllable and applicable actions of the robot, a multidimensional reward function and a Q -value model are proposed.
- We compare our RL method with others by evaluating them on a dataset of 1000 balls in simulation. An efficient retraining step is used to close the sim-to-real gap. Our models trained in real robots (see Fig. 1 right) achieve remarkable performance.

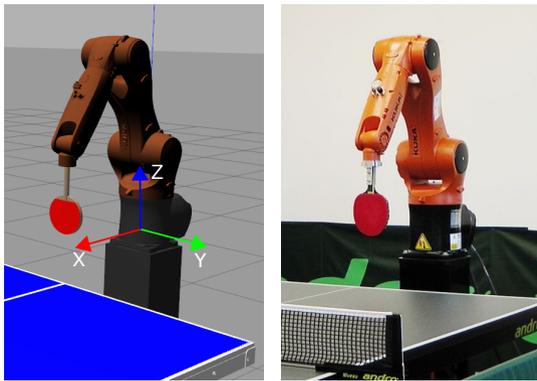


Fig. 1 Left: the simulated table tennis robot in a realistic simulation environment. To approximate reality, Gaussian noise is applied to the 3D position of each ball in simulation. The world coordinate system is identical to that of the robot. Right: our table tennis robot with a KUKA KR 6 R900 robot.

2 Related work

2.1 Simulation for robotic table tennis

Just as simulation is being used in various areas of robotics research (e.g., autonomous vehicles [5], industrial robots [23], and humanoid robots [24]), it is also being widely used for facilitating the development of more intelligent robotic table tennis [12, 25]. By using simulation, the robot can explore the action space extensively without worrying about safety during the training steps. Moreover, simulation provides the ability to compare the performance of different approaches in a fair and deterministic environment. For instance, [13] proposed an approach to sample-efficient learning of complex policies in the

context of robotic table tennis. The simulation environment was created using PyBullet [26] and connected to a virtual reality setup to capture human actions with instrumented paddles. Similarly, a simulated robotic table tennis system was built using PyBullet [26] to train policies for performing table tennis ball return tasks [12]. In addition, a simulation environment was developed in MATLAB and used to generate the optimal trajectory for robot table tennis [27]. To make use of robotic drivers and devices, [28] developed a flying robot with the Gazebo simulator which can be easily combined with ROS. [16] presented a hybrid simulation and real (HYSR) training method for muscular robots performing a table tennis task. To teach robots how to play table tennis without using real balls, the historical states of the ball in reality were recorded and replayed in simulation, and then the actions in simulation were applied to real robots. The aforementioned works are further evidence of the benefit and importance of using simulations in robotic table tennis research. However, some of these existing simulated table tennis systems have the characteristic that a high effort is required when transferring the trained models from simulation to reality due to the fact that they are not sufficiently realistic. In addition, due to incompatibility with the existing RL libraries that contain advanced RL algorithms, some of these existing simulated table tennis systems are further limited in their generalizability. Therefore, there is an increasing need to develop a more realistic simulation that can better simulate the real situation for advanced robot table tennis research. In this case, we developed a realistic simulation in combination with the Gazebo, ROS, OpenAI Gym and Spinning Up. The dynamics of the ball in simulation is determined based on [29, 30]. The simulated manipulator in our work is controlled in Cartesian space so that it can be easily replaced by other types of robots, such as flying robots or mobile robots.

2.2 Reinforcement learning in robotic table tennis

Reinforcement learning has been shown to be an excellent method for training robots to learn complex tasks [31, 32], and therefore could be a solution for robotic table tennis training. Indeed, deep RL has already attracted great interest from researchers and has achieved some success in the field of robotic table tennis [33, 34]. For instance, [12] developed an end-to-end RL algorithm for learning efficient policies to directly control of a simulated table tennis robot in joint space. A multi-modal model-free policy was trained to learn the velocities of each joint at 100 Hz by taking the joint position trajectories and locations of the ball as inputs. The optimal policy was found at about 1 million episodes. In [16], a muscular table tennis robot was trained in joint space in a hybrid simulation and reality (HYSR) system. The PPO was used as backbone. To return and smash the ball with a high success rate, a dense reward function was developed depending on the ball position and the robot state. Without using real balls, the robots were trained to play table tennis. [13] incorporated stroke learning into a hierarchical control system using an inverse landing model, an analytic racket controller, a forward racket model

and a forward landing model. Each model was trained separately to make the learning process easier and more efficient. A striking policy that can hit the ball to the targets with reasonable errors was learned from about 7,000 demonstrated trajectories, and in addition, the agent can learn from about 24,000 strikes in self-play to make optimal use of the human dynamics models for longer play. In addition, to efficiently learn the optimal stroke, a two-stage approach was adopted in [14]. In the first stage, the hitting states of the ball (position and velocity) were determined by an extended Kalman filter (EKF) based predictor. In the second stage, these determined states of the ball were fed as inputs to the DDPG, with the velocities of the racket being the outputs. A reliable performance was achieved within 200,000 simulated episodes. Instead of using DDPG directly, [35] proposed an accelerated parametrized-reward gradients approach to learn the velocity and rotation of the racket from the predicted hitting states. The policy was trained with 200 human demonstrations. To keep the explorations safe and avoid collisions, the racket action were restricted to a narrow range, resulting in the network being trained using a set of very similar trajectories of the ball.

3 Methodology

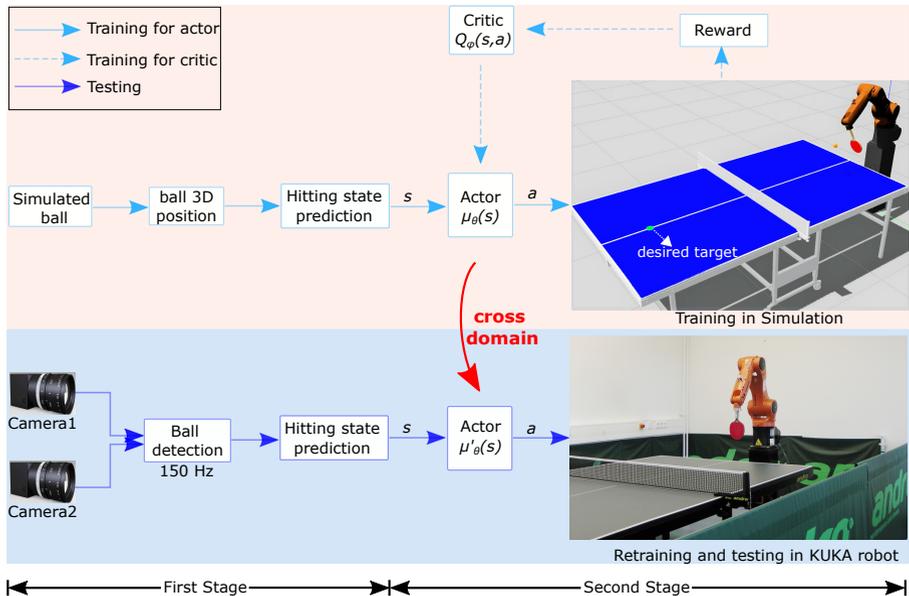


Fig. 2 The entire framework for training and testing. The goal is returning the ball to the desired target position on the table. In the first stage, the ball's state s at the hitting time is predicted by [30, 36]. In this work, we focus on the second stage, where an optimal stroke can be learned based on a novel RL algorithm. The upper part is performed in simulation. The RL model is trained within 10,000 episodes. To cross the reality gap, the model is then retrained in reality (see the bottom part).

To efficiently learn the optimal stroke and successfully return the ball to the desired target position on the table, we propose a novel framework as shown in Fig. 2. Specifically, a realistic simulation environment is developed for robot learning and comparison with other advanced RL algorithms. In the first stage, the hitting state (position, velocity, and spin) of the ball is predicted using the physical model of the ball proposed in our previous work [30, 36]. In the second stage, a novel approach is presented to learn the optimal stroke in simulation, which is conjugated with ROS and OpenAI libraries (see Fig. 3).

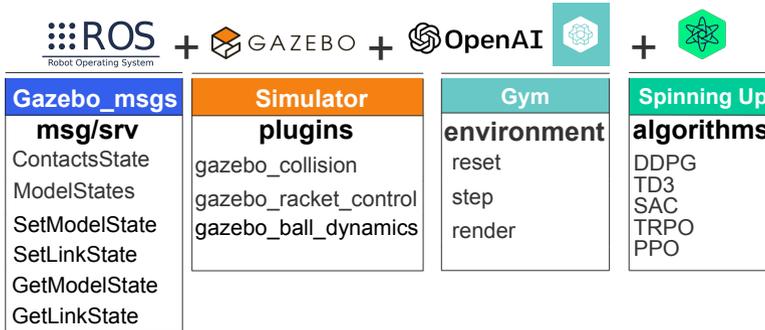


Fig. 3 The learning architecture is built by combining the Gazebo simulator with ROS, Gym, and Spinning Up. We can subscribe and publish the states (pose and velocity) of the ball and the robot from Gazebo with the `Gazebo_msgs` in ROS. These states should be fed into the Gym toolkit, which allows us to use the RL algorithms in Spinning Up.

3.1 Simulation

A well-known challenge for deep RL is the safe interaction with the environment. In particular, in robotic table tennis, it is difficult to explore all possibilities in reality, since unexpected collisions would destroy the mechanical robot parts. Moreover, the robot must interact with the environment over a large number of steps to learn a high-level policy. Taken together, this makes the application of deep RL in robotic table tennis more challenging. To address these issues, we develop a realistic simulation that provides a convenient scenario for optimal stroke learning as well as a comparison of different algorithms. The pose and velocity of the simulated racket are controlled by the Gazebo plugins. The dynamic model for each entity is obtained with the methods described in the following subsections.

3.1.1 Flying ball model

In addition to the gravitational force F_g , a flying ball is usually influenced by the Magnus force F_m and the air drag F_d [37]. As shown in Fig. 4, F_m is perpendicular to the spin axis and the flight direction, while F_d is opposite to the flight direction. These forces can be computed using the following formulas:

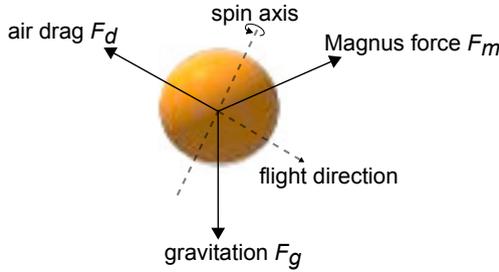


Fig. 4 Force analysis of a flying ball. A sphere shell of radius r_1 and mass m , with a centered spherical cavity of radius r_2 , is created as a simulated ball in Gazebo.

$$F_g = (0, 0, -mg)^T, \quad (1)$$

$$F_d = -\frac{1}{2}C_D\rho_a A\|v\|v, \quad (2)$$

$$F_m = \frac{1}{2}C_M\rho_a Ar_1(\omega \times v), \quad (3)$$

where the constants were determined in our previous work [36], including the mass of the ball $m = 2.7\text{g}$, the gravitational constant $g = 9.81\text{m/s}^2$, the drag coefficient $C_D = 0.4$, the density of the air $\rho_a = 1.29\text{kg/m}^3$, the lift coefficient $C_M = 0.6$, the radius of the ball $r_1 = 20\text{mm}$, the cavity radius $r_2 = 19.6\text{mm}$, and the cross-section of the ball $A = r_1^2\pi$. In addition, ω and v are the linear and angular velocities, respectively, which can be derived from the trajectory of the ball in reality. These two velocities can be further used to predict the hitting state of the ball using the algorithm introduced in [30, 36]. To simulate the accurate dynamics of the ball, the inertia value I should also be considered, which is calculated as

$$I = \frac{2}{5}m \left(\frac{r_1^5 - r_2^5}{r_1^3 - r_2^3} \right). \quad (4)$$

3.1.2 Bounce model

Since the physical contact between two objects in reality is a very complex matter, the Open Dynamics Engine (ODE), a popular rigid body dynamics library for robotics, was used to simulate the contact forces between the ball and the table (or the racket) in simulation. ODE has been originally integrated into the Gazebo simulator. To represent the elastic and frictional impacts on the ball, we compute the restitution coefficient κ_R and the friction coefficient μ similar to [38].

The restitution coefficient κ_R is defined as the ratio of the energy before and after a collision, e.g., when the ball bounces off the table. Approximately,

it can be solved by the free fall of the ball as follows:

$$\kappa_R^t = \frac{v_2^t - v_2^b}{v_1^b - v_1^t} = \frac{-v_2^b}{v_1^b} = \frac{-\sqrt{2 \cdot g \cdot h_2}}{-\sqrt{2 \cdot g \cdot h_1}} = \sqrt{\frac{h_2}{h_1}}, \quad (5)$$

where v_1^b and v_1^t are the velocities of the ball and table before bouncing, while v_2^b and v_2^t are the velocities after bouncing. h_1 and h_2 are the corresponding heights when the ball is not moving. Here the table velocity $v^t = 0$.

The friction coefficient μ is obtained from the setup in Fig. 5, where three balls are arranged together in the shape of a triangle frame. We first place the balls on the table and then lift the table until the balls start to slide. After obtaining the horizontal angle change θ of the table, the friction coefficient μ^t between the table and the ball can be calculated as

$$\mu^t = \frac{3mg \cdot \sin \theta}{3mg \cdot \cos \theta} = \tan \theta. \quad (6)$$



Fig. 5 Setup for measuring the friction coefficient μ . Three balls are arranged together in the shape of a triangle frame and placed on a flat table.

With the same methods, we calculate the restitution coefficient κ_R^r and the friction coefficient μ^r of the racket. The resulting parameter values are given in the Table 1. The additionally required parameters μ_2 and *slip* are defined as the friction coefficient in the second ODE friction pyramid direction and the coefficients of force-dependent-slip (FDS), respectively. They are manually adjusted to fit the reality.

Table 1 Collision parameters when the ball impacts on the table and racket.

	κ_R	μ	μ_2	<i>slip</i>
Table	0.97	0.05	0.025	0.01
Racket	0.9	1.0	0.025	0.01

To roughly test the accuracy of the simulation, we utilize a ball throwing machine to launch a topspin ball towards the stationary racket mounted on the robot. The entire trajectory of the ball can be recorded as the ground

truth using stereo cameras at 150 Hz. The starting spin and velocity of the ball are computed using a spin detector tool and a curve fitting approach, respectively. Then, these initial parameters are fed into the simulated ball and this initialized ball is served in simulation. The simulated trajectory is generated and is shown in Fig. 6.

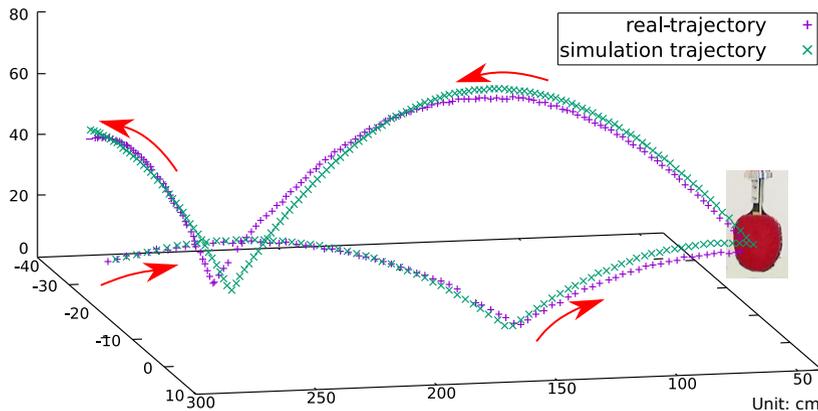


Fig. 6 Trajectory comparison between reality and simulation when serving a topspin ball. The difference between the returned landing positions on the table is about 6.2 cm. In this case, the racket is stationary since it is difficult to understand and simulate all the dynamics parameters of a moving racket in simulation. Therefore, the actions of the racket applied in simulation are always the true actions without any noise and time latency. This reality gap is closed in the retraining section 4.3.

3.2 Algorithm

With regard to the different types of inputs, there are generally two ways available when using deep RL algorithms in robotic table tennis. One is the one-stage algorithm that takes the state of the ball at each step as input and learns the pose of the racket in an end-to-end way. The other is the two-stage algorithm that first predicts the hitting state (i.e. position, velocity, and spin) of the ball and then takes it as input for RL. The latter can significantly accelerate the training step and can efficiently deal with different spin balls. In this work, we adopt the second approach to learn the optimal stroke of the racket based on the prediction of the state of the ball at the hitting time.

Since there is only a single state vector as input in the second stage, we then parameterize stroke learning as a bandit problem, where actions have no influence on the next states and consequently there are no delayed rewards in an episode. It is a simple version of an Markov Decision Process (MDP), with

$$M = (S, A, R), \quad (7)$$

where S is the set of the observed 11-D states s , including the 3D position p^b of the ball, 3D linear velocity v^b , 3D angular velocity ω^b at the hitting time,

and the desired 2D landing target position p^{tar} on the table. A is the set of 3D actions a that can be performed on the robot. Due to the restriction of the current mechanical structure and the control system, we cannot operate the robot as flexibly as a human can move. Therefore, we only learn to change the linear velocity v_x^r of the robot along the x -axis and its orientation angles (β^r, γ^r) around the y and z axes. The target position of the racket is the same as the predicted hitting position of the ball. R is a set of the immediate rewards r .

We use a policy $\mu_\theta(s)$ as the actor network which can output the actions a with respect to the current state s , as shown in Fig. 7 left. To evaluate the actions, a critic network $Q_\phi(s, a)$ is used, which takes as input both the states and actions and outputs a Q -value vector, as shown in Fig. 7 right. θ and ϕ are the weights of the neural network. The goal of our work is to learn a deterministic policy $\mu_\theta(s)$, which provides an action that maximizes the norm of $Q_\phi(s, a)$. According to the DDPG algorithm, the critic and the actor can be updated, respectively, by minimizing the losses:

$$\mathcal{L}(\phi, \mathcal{B}) = \mathbb{E}_{(s,a,r) \sim \mathcal{B}} \|Q_\phi(s, a) - r\|^2, \quad (8)$$

$$\mathcal{L}(\theta, \mathcal{B}) = -\mathbb{E}_{s \sim \mathcal{B}} \|Q_\phi(s, \mu_\theta(s))\|, \quad (9)$$

where B is a minibatch for storing s, a, r . The reward r is the feedback from the environment, which we will discuss in more detail later.

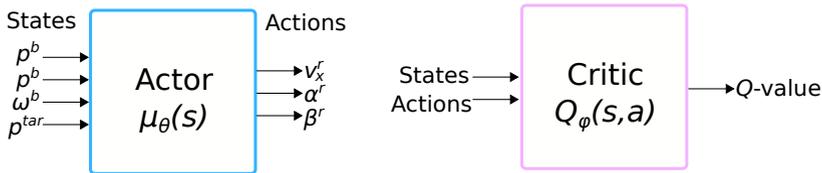


Fig. 7 Classic Actor-Critic algorithms. Instead of a 1D Q-value, a 3D Q-value is proposed to train the corresponding 3D actions.

To accelerate the training step and boost the resulting performance, we apply the following modifications to the classic actor-critic algorithms for training.

3.2.1 Exploration

For continuous action spaces, several exploration strategies are used in the deterministic environments. For instance, the *random* strategy selects actions randomly from a Gaussian distribution, while the *epsilon - greedy* strategy takes the random actions occasionally with probability ϵ and uses the output

of the current actor $\mu_\theta(s)$ with probability $1 - \epsilon$. This is the default strategy used in Spinning Up. In this work, we noticed that the actor $\mu_\theta(s)$ did not give the action with the maximum Q -value in the earlier training step because of the large loss error. Therefore, we generate the action a as

$$a = \operatorname{argmax}_{\mu_\theta(s) + \mathcal{N}} \|Q_\phi(s, \mu_\theta(s) + \mathcal{N})\|, \quad (10)$$

where \mathcal{N} is a Gaussian noise.

3.2.2 Reward shaping

In [14], a reward function was developed that depends on the height h^b of the ball when crossing the net and the actual landing position p^{real} on the table when the ball is returned. To balance these two independent variables h^b and p^{real} , a coefficient is used; however, deciding the value of this coefficient is a tricky issue. To address this problem, we decompose the reward function into three vector components: r_x and r_y for the x and y of the landing position, and r_h for the height of the ball. Each reward function is then normalized to $[0,1]$ by the following equations:

$$r_x = e^{-|p_x^{real} - p_x^{tar}|}, \quad (11)$$

$$r_y = e^{-|p_y^{real} - p_y^{tar}|}, \quad (12)$$

$$r_h = e^{-|h^b - 0.173|}, \quad (13)$$

$$r = [r_x, r_y, r_h] \quad \text{if } success \text{ else } \vec{0}, \quad (14)$$

where p_x and p_y are the landing position (in meters) along the x and y axes, the constant 0.173 is the measured actual height of the net (in meters), and e is the natural exponential operation. When a ball is successfully returned to the opposite table, *success* is set to true.

3.2.3 3D Q -value

Normally, the Q -value is a 1D vector that is expected to be maximized. To take advantage of the above rewards, we replace the last layer in the critic network from 1D to 3D. This results in a 3D Q -value $[Q_x, Q_y, Q_h]$, which can precisely indicate the quality of the actions.

In addition, for the actor-critic model we adopt TD3 as the backbone. The critic is changed to:

$$Q_\phi(s, a) = \begin{cases} Q_{\phi_1}(s, a), & \text{if } \|Q_{\phi_1}(s, a)\| < \|Q_{\phi_2}(s, a)\| \\ Q_{\phi_2}(s, a), & \text{otherwise} \end{cases}. \quad (15)$$

The whole training process is depicted in Algorithm 1, where the loss functions $\mathcal{L}(\phi_i, \mathcal{B})$ and $\mathcal{L}(\theta, \mathcal{B})$ are used to update the critic and the actor, respectively.

Algorithm 1 Policy Gradient Training with TD3 backbone

Input: Initial actor weights θ , critic weights ϕ_1 and ϕ_2 , replay buffer \mathcal{D} , number of episodes λ , Gaussian noise $\mathcal{N}(0, 0.1)$

Output: Optimal policy $\mu_\theta^*(s)$

```

1: for n=1 to  $\lambda$  do
2:   Observe the state  $s$  and generate the action  $a$  by

3:     
$$\operatorname{argmax}_{\mu_\theta(s) + \mathcal{N}} \|Q_\phi(s, \mu_\theta(s) + \mathcal{N})\|$$


4:   Apply  $a$  in the environment and get the reward  $r$ 
5:   Store  $(s, a, r)$  in the replay buffer  $\mathcal{D}$ 
6:   Reset the environment
7:   if it is time to update then
8:     Sample a random minibatch  $\mathcal{B}$  from  $\mathcal{D}$ 
9:     for  $i=1, 2$  do
10:      Update the critic by minimizing the loss:

11:      
$$\mathcal{L}(\phi_i, \mathcal{B}) \& = \mathbb{E}_{(s, a, r) \sim \mathcal{B}} \left[ \|Q_{\phi_i}(s, a) - r\|^2 \right]$$


12:     end for
13:     if it is time to update actor then
14:      Update the actor by minimizing the loss:

15:      
$$\mathcal{L}(\theta, \mathcal{B}) \& = -\mathbb{E}_{s \sim \mathcal{B}} \|Q_\phi(s, \mu_\theta(s))\| ,$$


16:      where  $Q_\phi(s, \mu_\theta(s))$  is got from Eq. 15:
17:     end if
18:   end if
19: end for

```

4 Experiments

4.1 Training and testing

To generalize the trained model, 11,000 serves were randomly sampled from a wide range of values for training (10,000) and evaluation (1,000), as shown in Table 2. To bridge the reality gap, a Gaussian noise is first applied to the 3D position of each ball in simulation. Then, we predict the state of the ball at the hitting time using the methods in [30, 36] instead of using the actual state of the ball in simulation. Since the predicted hitting position is the one to which

the simulated racket should actually move, this allows the replication of the real situation and makes the trained model more realistic for the real world. The final state variables and their range of the ball are shown in Table 2, which includes the desired landing target position (p_x^{tar}, p_y^{tar}), the position of the ball (p_x^b, p_y^b, p_z^b), linear velocity (v_x^b, v_y^b, v_z^b) and angular velocity ($\omega_x^b, \omega_y^b, \omega_z^b$) at the hitting time. These state variables are then normalized and used as inputs for training. The p_x^b is consistently equal to 0.675 m to form a virtual hitting plane used for the state prediction in the first stage. Thus, the hitting time is the time when the ball reaches this virtual plane. The home position of the racket in the world coordinates is (0.5, 0.0, 0.0) in meters.

Table 2 State variables and their range at the hitting time for training and evaluation in simulation.

	training	evaluation
p_x^{tar}	2.55m	
p_y^{tar}	0.0m	
p_x^b	0.675m	
p_y^b	[-0.60m, 0.63m]	[-0.68m, 0.68m]
p_z^b	[-0.01m, 0.34m]	[-0.01m, 0.34m]
v_x^b	[-6.00m/s, -1.35m/s]	[-5.94m/s, -2.52m/s]
v_y^b	[-1.95m/s, 2.16m/s]	[-1.29m/s, 2.02m/s]
v_z^b	[-3.47m/s, 3.15m/s]	[-3.40m/s, 2.60m/s]
ω_x^b	[-127.67rad/s, 110.88rad/s]	[-95.08rad/s, 111.53rad/s]
ω_y^b	[-299.99rad/s, 299.81rad/s]	[-299.62rad/s, 299.73rad/s]
ω_z^b	[-193.81rad/s, 189.65rad/s]	[-189.05rad/s, 189.47rad/s]
Episodes λ	10000	1000

Considering the mechanical setup of the robot, we restrict the linear velocity v_x^r of the robot to a range from 0m/s to 2m/s. The orientation angles β^r, γ^r are between -50° to 50° . The third angle α^r around the x -axis is calculated as

$$\alpha^r = k \cdot \frac{p_y^b}{0.5 \cdot w^t}, \quad (16)$$

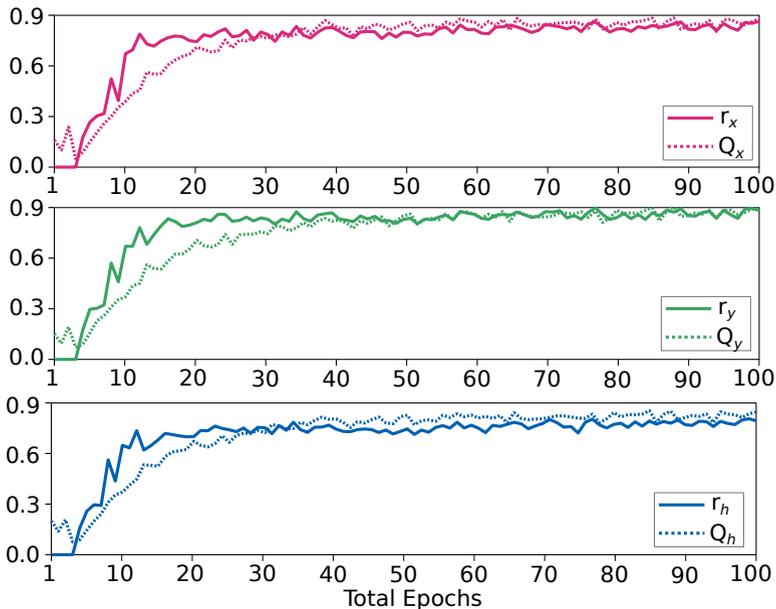
where w^t is the table width and k is a weight. In this way, the robot generates a flexible stroke. Here we assume the angle α^r has no influence on the impact with the ball.

In Equation 10, the added action noise \mathcal{N} for exploration is a mean-zero Gaussian distribution with a standard deviation of 0.1. The replay buffer \mathcal{D} has a size of 5,000. The number of training episodes λ is 10,000. Other hyper-parameters used for actor-critic are given in Table 3. The output actions from the actor are scaled to the valid range and then applied to the simulation. These hyper-parameters are tuned manually in order to achieve the best performance.

Compared to other environments that require millions of interactions in OpenAI Gym, our model is able to converge after 30 epochs, which took about

Table 3 Hyper-parameters used for training in simulation and retraining in reality.

	Actor/Critic	
	Training	Retraining
batch size	512	50
epochs	100	-
episodes per epoch	100	20
learning rate	1e-4	5e-5
optimizer	Adam	
layers	[256,256,3]	
activation	relu	
output activation	tanh/linear	

**Fig. 8** The testing rewards $[r_x, r_y, r_h]$ and the 3D Q -value $[Q_x, Q_y, Q_h]$ over the epochs in simulation.

1 hour of training. In addition, 1000 episodes were run for evaluation after each epoch. The resulting rewards and the corresponding 3D Q -value are plotted in Fig. 8. It is observed that the testing rewards reach a stable level starting from the 20th epoch, although the Q -values have not yet converged to the maximum values.

4.2 Evaluation

A commonly used metric for evaluation is the distance error between the actual and desired landing position [13, 14, 16]. However, this metric cannot reflect a failed return, e.g., if the landing position is not on the table, then it will be difficult to calculate the distance error. Therefore, we introduced a new metric

here: distance error ϵ_d computed by

$$r_d = e^{-\|p^{real}-p^{tar}\|_2} \quad \text{if success else } 0, \quad (17)$$

$$\epsilon_d = -\ln \left(\frac{1}{\lambda} \sum_{n=1}^{\lambda} r_d^n \right), \quad (18)$$

where the r_d^n is the reward for the landing distance error of the n^{th} ball. The value of the r_d^n is 0 if the ball does not land on the opposite table. Furthermore, two additional metrics reflecting the performance of the model were used for evaluation. The flying height error ϵ_h of the ball when crossing the net is calculated by

$$\epsilon_h = -\ln \left(\frac{1}{\lambda} \sum_{n=1}^{\lambda} r_h^n \right). \quad (19)$$

As well, the success rate of the ball in returning to the opposite table is calculated. To obtain a fair evaluation, we adopt 1000 episodes to cover a large range for the states (see Table 2 right).

Since only the 1D Q -value was used in the existing RL algorithms, we thus create a 1D reward function similar to [14] as

$$r_{eval} = e^{-k(\|p^{real}-p^{tar}\|-|h^b-0.173|)}, \quad (20)$$

where k is a scalar coefficient set to 0.5 in this work. This new reward function is used only for training the existing RL algorithms, including TRPO, PPO, SAC, DDPG, and TD3. By evaluating the existing and newly introduced algorithms using different reward functions corresponding to these different algorithms, we compute the distance error ϵ_D , the height error ϵ_h , and the success rate, respectively, as shown in Table 4. The unit of these errors is converted from meters to centimeters for better visualization. The proposed approach, argmax exploration plus 3D Q -value together with TD3 backbone, achieves better performance than the DDPG backbone. The other three approaches, TRPO, PPO, and SAC, learn the optimal stroke using a stochastic policy, resulting in much higher errors and lower success rate.

4.3 Retraining in reality

Although we built a high-fidelity simulation by manually measuring the coefficients and applying random noise to the ball, the real robot has many more dynamic and complicated factors that cannot be accurately measured and accounted for. To find the best hyper-parameters for retraining, we change the racket's restitution coefficient κ_R^r and friction coefficient μ^r in simulation. In this way, we can replicate the situation between two different rackets in reality. Based on the pretrained actor-critic model, we then retrain the model in

Table 4 Evaluation for different Algorithms.

Algorithms	ϵ_d	ϵ_h	success rate
TRPO	47.0cm	31.0cm	84.8%
PPO	44.2cm	30.8cm	87.1%
SAC	43.5cm	29.0cm	89.2%
DDPG	25.6cm	22.1cm	95.6%
DDPG+argmax	23.0cm	22.3cm	97.4%
DDPG+argmax+3D Q-value	21.3cm	21.7cm	97.9%
TD3	25.2cm	22.3cm	97.2%
TD3+argmax	22.2cm	21.2cm	97.7%
TD3+argmax+3D Q-value	20.3cm	21.2cm	98.5%

the new simulation with different batch sizes, episodes per epoch, and learning rates. The best hyper-parameters found in simulation are shown in Table 3 right.

A ball throwing machine, TTmatic 404A, is used to provide a variety of balls with sidespin, topspin, and backspin using a group of selected parameters. At the moment our robot can only handle sidespin and topspin, since the backspin ball causes too much acceleration in a robot joint. This could be solved in the future. In addition, the Reflexxes motion library [39] is used for robot trajectory planning in Cartesian space. Each epoch includes both sidespin and topspin balls during retraining. The state variables and the range of the ball for retraining and testing at the hitting time are shown in Table 5. Here, the model is retrained with 20 epochs in 0.5 hours to ensure that it achieves convergence. The hitting position p_x^b along the x axis is fixed to 0.675 m. The retraining process demonstrating the landing distance error ϵ_d and the height error ϵ_h is shown in Fig. 9.

Table 5 State variables and their range at the hitting time for retraining and testing in reality.

	retraining/testing in machine	testing with human
p_y^b	[-0.55m, 0.64m]	[-0.65m, 0.43m]
p_z^b	[0.085m, 0.34m]	[0.06m, 0.0.33m]
v_x^b	[-5.20m/s, -3.5m/s]	[-5.6m/s, -2.9m/s]
v_y^b	[-1.05m/s, 2.35m/s]	[-2.38m/s, 1.25m/s]
v_z^b	[-0.78m/s, 3.92m/s]	[-0.4m/s, 2.48m/s]
ω_x^b	[-32.94rad/s, 52.68rad/s]	[-33.00rad/s, 78.48rad/s]
ω_y^b	[-210.52rad/s, 5.33rad/s]	[-182.72rad/s, -55.28rad/s]
ω_z^b	[-157.65rad/s, 34.51rad/s]	[-66.68rad/s, 52.62rad/s]

Furthermore, to investigate the generalizability of the algorithms for a coefficient-unknown racket, we retrain a new model for a second racket whose dynamics are completely different from the first. Fig. 9 (b) illustrates the second retraining process. As can be seen, the retraining of the second racket required around 5 more epochs to converge and achieved similar performance to the retraining of the first racket.

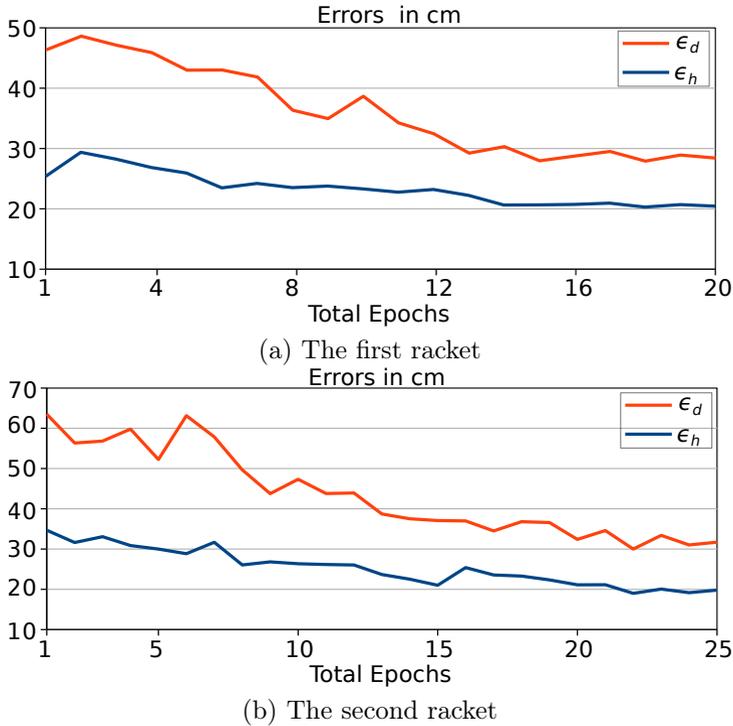


Fig. 9 Retraining process for the original racket (a) and for the coefficient-unknown racket (b), with balls served by a ball throwing machine. The metrics ϵ_d and ϵ_h in units of centimeters (cm) are the landing distance error ϵ_d and the height error ϵ_h when the ball crosses the net.

4.4 Testing in reality

To test our algorithms comprehensively, we conduct the experiments in three scenarios (see Fig. 10) in which the complexity gradually increases. In the



Fig. 10 Testing in three scenarios: Player1, Player2, and Machine.

first scenario, a human player (Player1) serves the ball with different starting positions at the front of the table. In this way, the hitting position along the y -axis can be completely covered. In the second scenario, a human player (Player2) with higher skills than Player1 plays a long game rally to test the continuous performance of the robot. The state variables and the range of the ball for these two scenarios are shown in the third column of the Table 5. In

the last and most complex scenario, the ball throwing machine (Machine) is used to serve balls with various spins and speeds. Since many aspects such as the robot, the racket, the state of the ball, the human players as well as the evaluation metrics are completely different, it is difficult to fairly compare the performance of our algorithm with other previous works. Therefore, we report the Table 6 by directly using the data in [40, 41] or by manually analyzing from [16, 35].

Table 6 Testing in reality.

Scenarios	episodes	ϵ_d	ϵ_h	success rate
Büchler et al. [16]	107	76.9cm	-	75%
Mülling et al. [40]	30	46.0cm	-	97%
Kyohei et al. [41]	100	22.5cm	-	99%
Jonas et al. [35]	300	24.2cm	-	93.6%
Player1	40	20.3cm	22.2cm	98%
Player2	40	25.6cm	23.5cm	
Machine	40	28.8cm	20.2cm	

The average distance error ϵ_d and height error ϵ_h in the testing in three real scenarios are 24.9 ± 9.0 cm and 21.9 ± 4.6 cm, respectively. Playing performance including some failure cases can be found at <https://youtu.be/SNnqtGLmX4Y>.

5 Conclusions and future work

In this work, we developed a realistic simulation for a table tennis robot. To learn the optimal stroke movement for the robot, we proposed a new policy gradient approach with TD3 backbone. Different algorithms were fairly evaluated in simulation using 1000 balls with a wide range of spins and speeds. To cross the domain from simulation to reality, a retraining approach was employed for the original racket and a coefficient-unknown racket racket. The test experiments showed a successful return rate of 98% in three complicated scenarios. The total training time is about 1.5 hours, which means that our algorithm is very efficient for application in robotic table tennis. Instead of a constant target position on the table, one can simply train a higher-level policy with a random target to make the game more challenging for human players. Moreover, our approach can be easily adapted to other robots playing racket-based sports, such as tennis, badminton, or squash.

Although we have shown significant improvements in robotic table tennis in both simulation and reality, the current control approach and the mechanical structure of the robot still limits its application in the real world. For example, the robot will fail to return the ball if the incoming ball is too high or too slow, since the target cannot be reached at a fast enough speed. Also, our robot will not have sufficient reaction time if the ball is too fast (e.g. $>10\text{m/s}$) because the minimum communication time between the controller and the robot is 5 ms. In the future we plan to optimize the Reflexxes motion libraries to produce

a more applicable trajectory for back spin balls. Instead of constraining the hitting position p_x^b along the x-axis, we will parameterize it as one action to be learned in RL. For a more offensive stroke, we could also try to learn the angular velocities of the racket so that the robot can initiatively generate a spin ball.

6 Acknowledge

We acknowledge the support of the Vector Stiftung and the KUKA Robotics Corporation.

References

- [1] Bertsekas, D.P., Tsitsiklis, J.N.: Neuro-dynamic programming: an overview. In: Proceedings of 1995 34th IEEE Conference on Decision and Control, vol. 1, pp. 560–564 (1995). IEEE
- [2] Sutton, R.S., Barto, A.G., *et al.*: Reinforcement learning. Journal of Cognitive Neuroscience **11**(1), 126–134 (1999)
- [3] Barto, A.G., Mahadevan, S.: Recent advances in hierarchical reinforcement learning. Discrete event dynamic systems **13**(1), 41–77 (2003)
- [4] Kendall, A., Hawke, J., Janz, D., Mazur, P., Reda, D., Allen, J.-M., Lam, V.-D., Bewley, A., Shah, A.: Learning to drive in a day. In: 2019 International Conference on Robotics and Automation (ICRA), pp. 8248–8254 (2019). IEEE
- [5] Osiński, B., Jakubowski, A., Zięcina, P., Miłoś, P., Galias, C., Homoceanu, S., Michalewski, H.: Simulation-based reinforcement learning for real-world autonomous driving. In: 2020 IEEE International Conference on Robotics and Automation (ICRA), pp. 6411–6418 (2020). IEEE
- [6] Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., *et al.*: Mastering the game of go without human knowledge. nature **550**(7676), 354–359 (2017)
- [7] Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C., *et al.*: Dota 2 with large scale deep reinforcement learning. arXiv preprint arXiv:1912.06680 (2019)
- [8] Gu, S., Holly, E., Lillicrap, T., Levine, S.: Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In: 2017 IEEE International Conference on Robotics and Automation (ICRA), pp. 3389–3396 (2017). IEEE

- [9] Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., *et al.*: Scalable deep reinforcement learning for vision-based robotic manipulation. In: Conference on Robot Learning, pp. 651–673 (2018). PMLR
- [10] Koos, S., Mouret, J.-B., Doncieux, S.: Crossing the reality gap in evolutionary robotics by promoting transferable controllers. In: Proceedings of the 12th Annual Conference on Genetic and Evolutionary Computation, pp. 119–126 (2010)
- [11] Cutler, M., How, J.P.: Efficient reinforcement learning for robots using informative simulated priors. In: 2015 IEEE International Conference on Robotics and Automation (ICRA), pp. 2605–2612 (2015). <https://doi.org/10.1109/ICRA.2015.7139550>
- [12] Gao, W., Graesser, L., Choromanski, K., Song, X., Lazic, N., Sanketi, P., Sindhvani, V., Jaitly, N.: Robotic table tennis with model-free reinforcement learning. In: 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 5556–5563 (2020). <https://doi.org/10.1109/IROS45743.2020.9341191>
- [13] Mahjourian, R., Miikkulainen, R., Lazic, N., Levine, S., Jaitly, N.: Hierarchical policy design for sample-efficient learning of robot table tennis through self-play. arXiv preprint arXiv:1811.12927 (2018)
- [14] Zhu, Y., Zhao, Y., Jin, L., Wu, J., Xiong, R.: Towards high level skill learning: Learn to return table tennis ball using monte-carlo based policy gradient method. In: 2018 IEEE International Conference on Real-time Computing and Robotics (RCAR), pp. 34–41 (2018). IEEE
- [15] Hanna, J.P., Desai, S., Karnan, H., Warnell, G., Stone, P.: Grounded action transformation for sim-to-real reinforcement learning. Machine Learning, 1–31 (2021)
- [16] Büchler, D., Guist, S., Calandra, R., Berenz, V., Schölkopf, B., Peters, J.: Learning to play table tennis from scratch using muscular robots. arXiv preprint arXiv:2006.05935 (2020)
- [17] Schulman, J., Levine, S., Abbeel, P., Jordan, M., Moritz, P.: Trust region policy optimization. In: Proceedings of The 32nd International Conference on Machine Learning, pp. 1889–1897 (2015)
- [18] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017)
- [19] Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement

- learning. arXiv preprint arXiv:1509.02971 (2015)
- [20] Fujimoto, S., Hoof, H., Meger, D.: Addressing function approximation error in actor-critic methods. In: International Conference on Machine Learning, pp. 1587–1596 (2018). PMLR
- [21] Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In: International Conference on Machine Learning, pp. 1861–1870 (2018). PMLR
- [22] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., Zaremba, W.: Openai gym. arXiv preprint arXiv:1606.01540 (2016)
- [23] Li, F., Jiang, Q., Zhang, S., Wei, M., Song, R.: Robot skill acquisition in assembly process using deep reinforcement learning. *Neurocomputing* **345**, 92–102 (2019). <https://doi.org/10.1016/j.neucom.2019.01.087>. Deep Learning for Intelligent Sensing, Decision-Making and Control
- [24] Abreu, M., Reis, L.P., Lau, N.: Learning to run faster in a humanoid robot soccer environment through reinforcement learning. In: Chalup, S., Niemueller, T., Suthakorn, J., Williams, M.-A. (eds.) *RoboCup 2019: Robot World Cup XXIII*, pp. 3–15. Springer, Cham (2019)
- [25] Gao, Y., Tebbe, J., Zell, A.: Robust stroke recognition via vision and imu in robotic table tennis. In: Farkaš, I., Masulli, P., Otte, S., Wermter, S. (eds.) *Artificial Neural Networks and Machine Learning – ICANN 2021*, pp. 379–390. Springer, Cham (2021)
- [26] Coumans, E., Bai, Y.: Pybullet, a python module for physics simulation in robotics, games and machine learning (2017)
- [27] Koç, O., Maeda, G., Peters, J.: Online optimal trajectory generation for robot table tennis. *Robotics and Autonomous Systems* **105**, 121–137 (2018)
- [28] Silva, R., Melo, F.S., Veloso, M.: Towards table tennis with a quadrotor autonomous learning robot and onboard vision. In: 2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 649–655 (2015). IEEE
- [29] Blank, P., Groh, B.H., Eskofier, B.M.: Ball speed and spin estimation in table tennis using a racket-mounted inertial sensor. In: Proceedings of the 2017 ACM International Symposium on Wearable Computers, pp. 2–9 (2017)

- [30] Tebbe, J., Klamt, L., Gao, Y., Zell, A.: Spin detection in robotic table tennis. In: 2020 IEEE International Conference on Robotics and Automation (ICRA), pp. 9694–9700 (2020). IEEE
- [31] Hester, T., Stone, P.: Texplora: real-time sample-efficient reinforcement learning for robots. *Machine learning* **90**(3), 385–429 (2013)
- [32] Gu, S., Holly, E., Lillicrap, T., Levine, S.: Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In: 2017 IEEE International Conference on Robotics and Automation (ICRA), pp. 3389–3396 (2017). <https://doi.org/10.1109/ICRA.2017.7989385>
- [33] Huang, Y., Büchler, D., Koç, O., Schölkopf, B., Peters, J.: Jointly learning trajectory generation and hitting point prediction in robot table tennis. In: 2016 IEEE-RAS 16th International Conference on Humanoid Robots (Humanoids), pp. 650–655 (2016). <https://doi.org/10.1109/HUMANOIDS.2016.7803343>
- [34] Yang, L., Zhang, H., Zhu, X., Sheng, X.: Ball motion control in the table tennis robot system using time-series deep reinforcement learning. *IEEE Access* **9**, 99816–99827 (2021). <https://doi.org/10.1109/ACCESS.2021.3093340>
- [35] Tebbe, J., Krauch, L., Gao, Y., Zell, A.: Sample-efficient reinforcement learning in robotic table tennis. arXiv preprint arXiv:2011.03275 (2020)
- [36] Tebbe, J., Gao, Y., Sastre-Rienietz, M., Zell, A.: A table tennis robot system using an industrial kuka robot arm. In: German Conference on Pattern Recognition, pp. 33–45 (2018). Springer
- [37] Zhang, Y., Zhao, Y., Xiong, R., Wang, Y., Wang, J., Chu, J.: Spin observation and trajectory prediction of a ping-pong ball. In: 2014 IEEE International Conference on Robotics and Automation (ICRA), pp. 4108–4114 (2014). IEEE
- [38] Blank, P., Groh, B.H., Eskofier, B.M.: Ball speed and spin estimation in table tennis using a racket-mounted inertial sensor. In: Proceedings of the 2017 ACM International Symposium on Wearable Computers, pp. 2–9 (2017)
- [39] Kröger, T.: On-Line Trajectory Generation in Robotic Systems: Basic Concepts for Instantaneous Reactions to Unforeseen (Sensor) Events vol. 58. Springer, ??? (2010)
- [40] Mülling, K., Kober, J., Kroemer, O., Peters, J.: Learning to select and striking movements in robot table tennis. *The International Journal of Robotics Research* **32**(3), 263–279 (2013)

- [41] ASAI, K., Nakayama, M., YASE, S.: The Ping Pong Robot to Return a Ball Precisely (2019). <https://www.omron.com/global/en/technology/omrontechnics/vol51/016.html>