

DexCatch: Learning to Catch Arbitrary Objects with Dexterous Hands

Fengbo Lan^{1,*}, Shengjie Wang^{1,2,3,*}, Yunzhe Zhang¹, Haotian Xu¹,
Oluwatosin Oseni⁴, Yang Gao^{1,2,3,†}, Tao Zhang^{1,†}, *Senior Member, IEEE*

Abstract—Achieving human-like dexterous manipulation remains a crucial area of research in robotics. Current research focuses on improving the success rate of pick-and-place tasks. Compared with pick-and-place, throw-catching behavior has the potential to increase picking speed without transporting objects to their destination. However, dynamic dexterous manipulation poses a major challenge for stable control due to a large number of dynamic contacts. In this paper, we propose a Stability-Constrained Reinforcement Learning (SCRL) algorithm to learn to catch diverse objects with dexterous hands. The SCRL algorithm outperforms baselines by a large margin, and the learned policies show strong zero-shot transfer performance on unseen objects. Remarkably, even though the object in a hand facing sideward is extremely unstable due to the lack of support from the palm, our method can still achieve a high level of success in the most challenging task. Video demonstrations of learned behaviors and the code can be found on the [supplementary website](#).

I. INTRODUCTION

The study of dexterous hands is a widely pursued research area aimed at matching the intricate capabilities of human hands. Researchers aspire to employ these hands for delicate manipulations [1]–[4]. In contrast to traditional planning approaches [5], reinforcement learning has witnessed remarkable success in robot planning tasks [6]–[9]. Consequently, researchers have applied reinforcement learning techniques to various tasks involving dexterous hands, leading to diverse outcomes [10]–[12]. Previous work has tackled a series of manipulation tasks, including in-hand reorientation [11], [13], opening a door [10], assembling LEGO blocks [14], and even solving a Rubik’s cube [12]. These tasks can be summarized as common behaviors in daily living, where hands pick up an object, reorient it in hand, and then either place it or use it as a tool [13].

While pick-and-place can fulfill the requirements of many daily tasks, dynamic throw-catching holds the potential for more efficient task execution [15]–[17]. For example, two robots can employ throwing and catching to efficiently transfer objects instead of pick-and-place, as shown in Fig. 1. Meanwhile, it can extend the working space of robots, which avoids the potential collisions for robot interactions [15]. However, precisely catching arbitrary objects proves to be highly challenging [18]. First, the dynamic contact with strong force makes the object fall from the hand. In previous work [18], model-free reinforcement learning

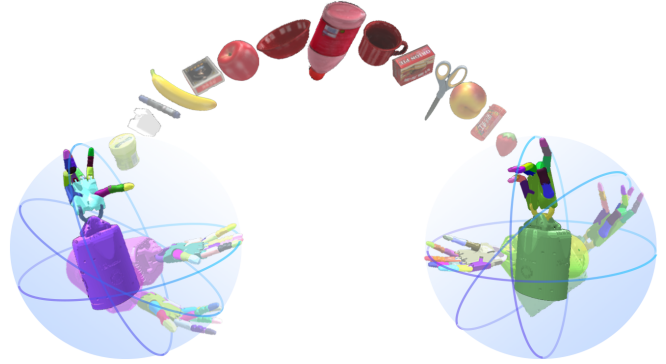


Fig. 1. Generalization tasks for throwing and catching objects with shadow hands. Shadow hands can throw and catch a variety of objects with different shapes, masses and other properties when the hands face sideward.

methods achieved nearly 0% success rate in the throwing-catching task within 30 million environment steps. Secondly, varying object-centric properties (e.g. mass distribution, friction, shape) require the catching behavior robust enough, especially for objects of daily life. Thirdly, faced with flying objects from random directions, the catching hand should use different postures. But the challenge is the task of catching objects with hands facing sideward, where objects easily fall from the hands due to the lack of any support surface, as illustrated in Fig. 1.

To address the challenge of dexterous two-handed throwing and catching, we present the Stability-Constrained Reinforcement Learning (SCRL) algorithm. This algorithm introduces three enhancements that significantly boost the exploration of stable catching behavior. 1) Stability Condition: The algorithm enforces stability in the state during exploration, enhancing the strategy’s robustness. 2) Intrinsic Advantage: Strategies are encouraged to explore actions with high advantage, contributing to more effective exploration. 3) Normalization Method: By constraining parameter representation space, the learning process is accelerated. Our method outperforms baselines on various throwing and catching scenarios where the hands face upward and sideward. Furthermore, the learned policy generalizes to create diverse catching gestures on the objects with varying properties (e.g., mass distribution, friction, shape). It is worth noting that it achieves a success rate for unseen objects, even on the most challenging catching task with the hands facing sideward.

II. RELATED WORK

The versatile dexterous hand plays a crucial role in performing various tasks across different human-focused set-

* Equal contribution; †Corresponding authors.

¹ Tsinghua University, Beijing, China; ²Shanghai Artificial Intelligence Laboratory, Shanghai, China; ³Shanghai Qi Zhi Institute, Shanghai, China;

⁴ Covenant University, Ogun State, Nigeria.

tings. However, effectively managing dexterous hand systems poses a challenge due to their complex actions and intricate contact interactions [5]. Traditionally, controllers for manipulation tasks heavily rely on dynamic models and trajectory optimization techniques [1]–[4]. This involves employing simplified dynamic models to generate control patterns using trajectory planning approaches. For instance, Williams *et al.* [19] achieved success using the Model Prediction Path Integration Control (MPPI) method to skillfully manipulate a cube. Charles Voss *et al.* [18] enhanced the MPPI method, making it simpler to manage task handovers between two hands.

In recent times, there has been a notable surge in the adoption of controller designs rooted in reinforcement learning. These approaches simplify the controller design process, and model-agnostic techniques have gained substantial popularity within the control domain. In the realm of dexterous manipulation, numerous endeavours have showcased remarkable advancements compared to conventional controllers [20]–[22]. Entities like OpenAI have harnessed reinforcement learning-based controllers to successfully rearrange blocks and solve Rubik’s cubes [12]. To expedite training, a blend of reinforcement learning and imitation learning has been employed to facilitate dexterous hands in acquiring diverse skills such as repositioning objects, utilizing tools, opening doors, and more [10]. Acknowledging the existing method’s limited ability to generalize, Chen *et al.* [13] introduced an effective system capable of learning how to reorient a wide array of objects. For the cooperation between two hands, Zakka *et al.* [23] and Chen *et al.* [22] proposed a series of challenging tasks for bi-manual dexterity, like mastering the piano and lifting the pot. Furthermore, a recent work introduced the Population-Based Training (PBT) algorithm in RL learning, leading to solving the two-handed reorientation task [21]. However, the work only focuses on a single object. Our method successfully achieves the throwing and catching of diverse objects, showcasing superior performance compared to baseline approaches.

III. DEXTEROUS CATCHING TASKS

A. Preliminaries

The throwing and catching tasks of dexterous hands can be formalized as an infinite horizon Markov Decision Process (MDP), which is defined by a tuple (S, A, R, P, ρ_0) with the transition function $P(\cdot|s, a)$. Among them, $S \in R^n$ is the current state space, $A \in R^m$ is the current action space, R is the reward, $\gamma \in [0, 1)$ is the discounted factor, and ρ is the initial state distribution. S' and A' is the next state space and the next action space. The goal of RL is to find an optimal cumulative return R of the MDP with the policy $\pi : S \leftarrow A$. The problem is shown to be represented as an optimization problem of the expected discounted cumulative return.

$$\mathbb{E}_{s_0 \sim \rho_0, a_t \sim \pi(\cdot|s_t), s_{t+1} \sim P(\cdot|s_t, a_t)} \left[\sum_{i=0}^{\infty} \gamma^i r(s_t, a_t) \right] \quad (1)$$

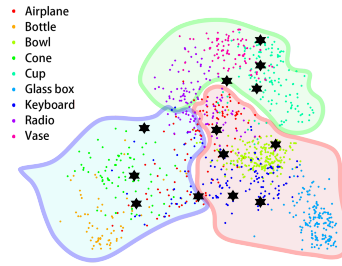


Fig. 2. The features of different point clouds are extracted by PointClipV2 [24], then reduced by PCA, and finally visualized in two dimensions. The Clustering Using Representatives(CURE) method shows that any object is approximately divided into three categories according to different shapes, extending in three directions, respectively, elongated, flat, and regular. Sample of objects each category are bottles, keyboards, and cups, respectively.

B. Task Description

Our research centers on the challenges of throwing and catching tasks using dexterous shadow hands. Our aim is for these shadow hands to adeptly handle diverse objects in various positions and execute agile catches using different gestures. To address this objective, our design comprises of three key components.

a) Environment Setting: The shadow hand models, which emulate human skeletons, possess 24 degrees of freedom (DOFs) and are operated using 20 pairs of agonist-antagonist tendons [13]. Meanwhile, the remaining four joints are underactuated. The presence of force sensors on the ten fingertips facilitates the acquisition of force-related data during catching and throwing maneuvers. In our paper, we investigate three catching scenarios, namely Abreast Catch, Underarm Catch and Overarm Catch, which are shown in Fig. 5 respectively.

In Abreast/Underarm Catch, a hand facing upward should throw an object to the other hand, and the other one should hold it in the palm later. Notably, Overarm Catch also solves the same problem but represents the most challenging catching task. This is because the object in the hands is not stable without the support from the palm when the hands are oriented vertically.

b) Object-generalized Setting: Our design objective mainly involves attaining the capability to generalize across diverse objects through the introduction of point cloud features of objects. To achieve this, we leverage the IsaacGym simulator [25], which enables parallel training, allowing the system to learn the skill of catching various objects simultaneously.

The 3D point cloud feature of an object represents its shape, created by densely gathered points in 3D space. We adopt point clouds to approximate object features. Our approach involves feature extraction via a point clouds processing network utilizing PointClipV2 [24]. This algorithm facilitates zero-shot classification segmentation to identify features. To reduce computational complexity, the extracted feature map degrades to two dimensions using Principal Component Analysis (PCA). Our experimentation and visual

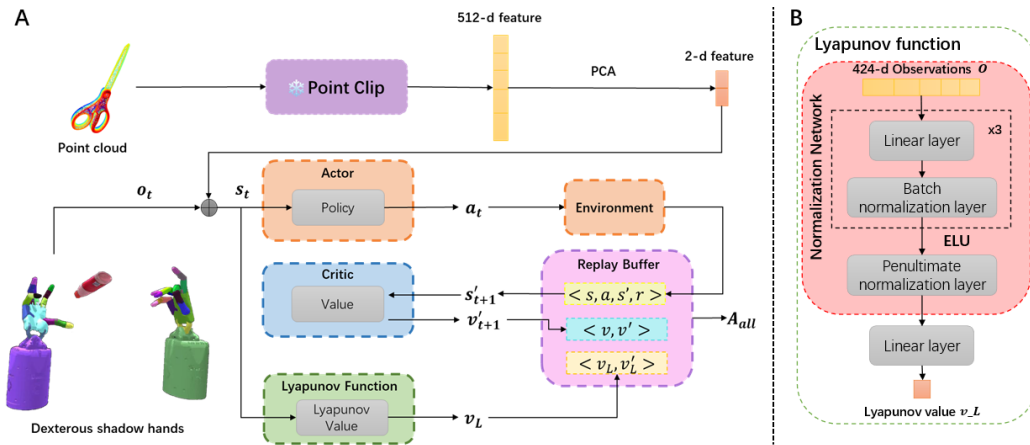


Fig. 3. **A:** The SCRL algorithm takes as input the environmental observation and the point cloud feature of the object. Then, it learns the catching policy for dexterous hands through an Actor-Critic structure. **B:** The Lyapunov function, the policy function and the value function are estimated using neural networks. Their network structures are similar, and the only difference is the dimension of the output. For example, the Lyapunov function contains a linear layer, a batch normalization layer, and a penultimate normalization layer, and the last linear layer takes a scalar as output.

TABLE I
REWARD DESIGN

r_p	$2 \arcsin\left(\frac{\ p_g - p_c\ _2}{\ q_d\ _2} - \frac{\ q_d\ _{min}}{\ q_d\ _{max}}\right) - \sum a_c^2$
r_o	
r_a	
r_{succ}	
r_{fall}	
r_{total}	$r_p + r_o + r_a + r_{succ} - r_{fall}$

analysis have substantiated that these 2-dimensional features adequately capture the distinct attributes of various objects, as demonstrated in Fig.2. It can be seen that the distribution of various objects extends along three distinct directions. The phenomenon allows for the classification of objects into three categories. The first category comprises elongated objects, exemplified by items like bottles. The second category encompasses flat objects, such as keyboards. Lastly, regular objects, which do not fit into the previous categories, such as cups, are classified under the third category. Furthermore, we add variations in the mass and inertia of objects during training, leading to the generalization ability about the different physical properties of objects.

c) State, Action and Reward: **1) State:** The state space consists of 424 elements. Among these, the 422-dimensional state incorporates crucial information such as hand position, hand orientation, angular and linear velocities of the hands, position and speed of each hand’s degree of freedom, force-torque sensor readings of the hands, target object’s position and orientation, as well as the current object’s position and orientation. The above observation corresponds to the setting in [22]. Additionally, a two-dimensional state features the input derived from the point cloud data of the object. It serves as the perceptual part for the agent regarding the object. **2) Action:** The action space comprises 52 dimensions, encompassing 20 pairs of agonist-antagonist tendons for both hands, along with three dimensions each for the required force and torque of each hand. **3) Reward:** The reward

is composed of five elements: 1) r_p quantifies the distance between the position of the goal object p_g and the current object p_c . 2) r_o represents the angular discrepancy between the target object’s orientation and the current object’s orientation, utilizing quaternions for accurate angular measurement. 3) r_a evaluates the smoothness of joint movement through the current actions a_c . 4) A constant reward is given for successfully completing tasks r_{succ} . 5) A penalty is imposed for task failure r_{fall} . Specifically, the smoothness of joint movement is controlled by negatively rewarding the sum of squares of action quantities as per the strategy. The overall reward is expressed r_{total} in Table I, where target position distance is calculated using Euler distance, angular distance is determined using quaternions, joint smoothness is penalized using action quantities’ sum of squares, and success and failure rewards are assigned as positive and negative constants, as outlined in Table I.

IV. STABILITY-CONSTRAINED REINFORCEMENT LEARNING ALGORITHM

To address the task outlined in the above section, we propose the Stability-Constrained Reinforcement Learning (SCRL) algorithm grounded in the Proximal Policy Optimization (PPO) [26] framework. Our approach integrates considerations for both system stability during task completion and efficient learning. Furthermore, we incorporate parameter space normalization to speed up the learning.

The algorithm’s architecture is illustrated in Fig. 3. Here’s a breakdown of its components: 1) Point Cloud Feature Extraction: The initial step involves extracting point cloud features using the PointClipV2 [24]. Subsequently, Principal Component Analysis (PCA) is employed to reduce dimensionality, yielding a 2-dimensional feature output. This output serves as the perception state for the shadow hands system; 2) Observation Input: The agent takes both the current observation of objects and the shadow hands and the compressed features from point clouds as input; 3) Actor-

Critic Framework: The actor and critic are represented as two neural networks, which is the same as PPO algorithm. The actor takes the current state s_t as input and executes the action a_t . The system interacts with the environment to obtain the reward r_t and the next candidate state s'_{t+1} . The critic evaluates the candidate state s'_{t+1} to estimate its value v' . Meanwhile, a Lyapunov function takes as input s_t , yielding the Lyapunov value v_L for the current state; 4) Replay Buffer and Parameter Update: Relevant data from the current state is stored in the replay buffer. The total advantage A_{all} is calculated based on r_t, v, v', v_L and v'_L , which plays a role in updating the parameters of the actor. The specific formulation of A_{all} is denoted in Eq. (10). The synergy of these steps constitutes the SCRL algorithm, achieving both stability and efficient exploration during task execution.

A. Proximal Policy Optimization(PPO)

The PPO algorithm is an on-policy reinforcement learning approach. It excels in robotics tasks featuring high-dimensional action spaces. The algorithm contains an actor and a critic network, where θ and ϕ denote the parameters of the actor network $\pi_\theta(a_t|s_t)$ and the critic network $V_\phi(s_t)$, respectively. For the policy function $\pi_\theta(a_t|s_t)$, the output corresponds to the mean and std of Gaussian distribution. On the other hand, the value function $V_\phi(s_t)$ yields the current state's value. Building upon the policy gradient method, the parameter updating formula for the actor network takes the following form:

$$\theta_{k+1} = \arg \max_{\theta} E_{\mathbf{s}, \mathbf{a} \sim \pi_{\theta_k}} [L(\mathbf{s}, \mathbf{a}, \theta_k, \theta)] \quad (2)$$

$$L(\cdot) = \min(\tilde{A}(\mathbf{s}, \mathbf{a}), \text{clip}(1, 1 - \epsilon, 1 + \epsilon)\tilde{A}(\mathbf{s}, \mathbf{a})) \quad (3)$$

$$\tilde{A}(\mathbf{s}, \mathbf{a}) = \frac{\pi_\theta(\mathbf{a}|\mathbf{s})}{\pi_{\theta_k}(\mathbf{a}|\mathbf{s})} A^{\pi_{\theta_k}}(\mathbf{s}, \mathbf{a}) \quad (4)$$

where $A^{\pi_{\theta_k}}(\mathbf{s}, \mathbf{a}) = G_t^{\pi_{\theta_k}} - V_\phi^{\pi_{\theta_k}}(s_t)$, is the advantage function, represents the tolerant range of deviation between the original and existing strategy. For parameters ϕ , the update is expressed as:

$$\phi_{k+1} = \arg \max_{\phi} E_{\mathbf{s}, \mathbf{a} \sim \pi_{\theta_k}} [G_t^{\pi_{\theta_k}} - V_\phi^{\pi_{\theta_k}}(s_t)] \quad (5)$$

However, because the calculation of $G_t^{\pi_{\theta_k}}$ value needs the whole data of a finished episode, we approximate $G_t^{\pi_{\theta_k}}$ to $r_t + V_\phi^{\pi_{\theta_k}}(s_{t+1})$ through TD learning method.

B. Lyapunov Stability

The Lyapunov function serves as a potent tool for assessing system stability. In the context of RL, the Markov Decision Process (MDP) represents a discrete-time dynamical system corresponding to continuous states, and the Lyapunov stability condition becomes a constraint within the objective function. This integration ensures both stability and optimality of the system. Indeed, it's important to note that in our framework, the points of optimality and stability do not always coincide. As depicted in Fig. 4, the middle-left figure illustrates a gesture that is optimal for the task but

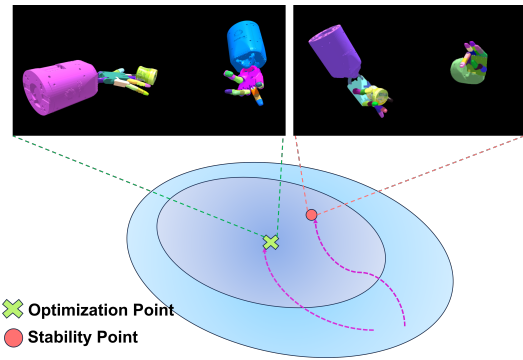


Fig. 4. The solution with stability and the optimal solution are not consistent in the task. For example, the optimal solution on the left, the reward and advantage remain the highest but the gesture is not stable. In the stable solution on the right, the hand totally hold the object.

lacks stability, whereas the right figure showcases a gesture that combines both stability and optimality.

One approach, known as the self-learning Lyapunov stabilization method [27], is utilized in our paper. First, we introduce a Lyapunov stability critic network, which represents the candidate Lyapunov function V_L using neural networks. Subsequently, the empirical Lyapunov risk $\mathcal{R}(\Lambda)$ [27] is defined as per Eq. (6). The parameter update is then executed by minimizing the loss function associated with the Lyapunov risk $\mathcal{R}(\Lambda)$, as part of achieving the desired system stability.

$$\mathcal{R}(\Lambda) = \frac{1}{N} \sum_{i=1}^N (\max(-V_L(s_i), 0) + \max(0, L_f V_L)) + V_L^2(0) \quad (6)$$

where $s_1, \dots, s_i, \dots, s_n$ is the state of sampling. The empirical Lyapunov risk is non-negative, we approximate the Lie derivative $L_f V_L$ by finite difference of the sampling trajectory of the system because of the complexity of modeling the dynamics of this system:

$$L_f V_L(s) = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} (V_L(s') - V_L(s)) \quad (7)$$

where s and s' are two consecutive states, and Δt is the time difference between them. Since the Lyapunov function is considered to be another value function, Lie derivative term $L_{f, \Delta t} V_L(s)$ is the only term in the Lyapunov risk affected by the policy. Therefore, an additional goal of ensuring $L_{f, \Delta t} V_L(s) < 0$ is added to the policy optimization, to encourage the policy's update with the requirement of stabilization. Hereafter, we design the part of the Lyapunov stability advantage function based on the clipped Lie derivative term:

$$A_L(s_t, a_t) = \min(0, -L_{f, \pi_\theta, \Delta t} V_L(s_t)) \quad (8)$$

C. Intrinsic Advantage

Based on the advantage function of the PPO algorithm, in order to improve the ability of the system efficiently, we design the intrinsic advantage function with reference to the intrinsic reward mechanism [28]. For a Markov decision

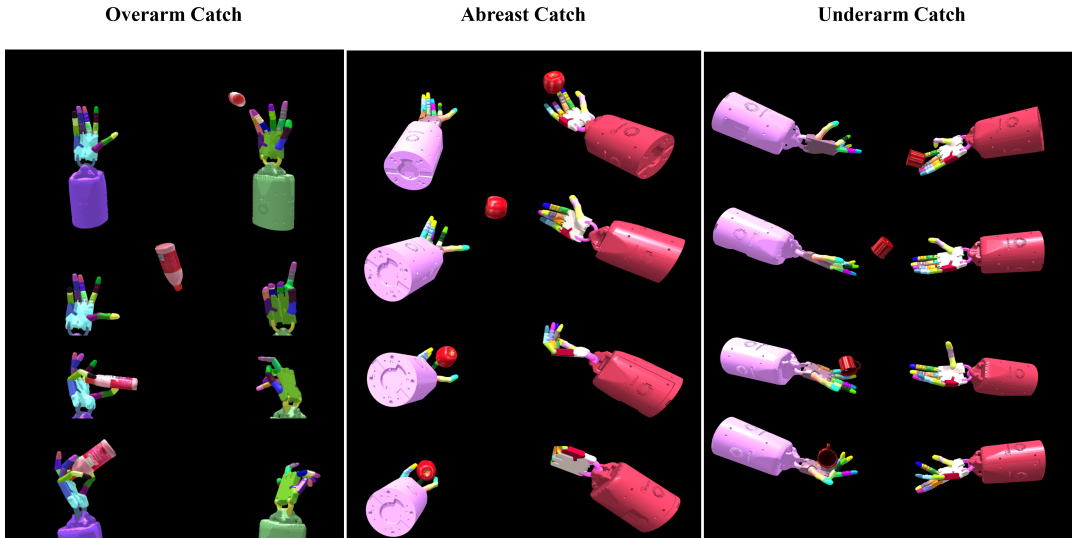


Fig. 5. Video snapshots of throwing and catching tasks performed by our method. There exist three tasks with different postures of hands, Overarm Catch, Abreast Catch and Underarm Catch. The most challenging task is Overarm Catch because the hands are oriented vertically. Notably, our method can catch diverse objects in three tasks.

process trajectory, different values $V_\phi(s_t)$ will affect the success of the final result of this trajectory. To provide a specific example, let's consider the fact that the next states s'_{t+1} exhibit varying conditions, leading to different values $V_\phi(s_{t+1})$. It's worth noting that in a successful process, the value $V_\phi(s_{t+1})$ tends to be higher. Consequently, within a Markov decision process, we can further refine the advantage function $A(s, a)$ using empirical measures. If an action is effective, it should result in a higher value, as represented by $V_\phi(s_{t+1}) - V_\phi(s_t) > 0$. This signifies that action a_t transitions to a state with a higher success rate. Conversely, an unfavorable action can lead to a decrease in value, reflected as $V_\phi(s_{t+1}) - V_\phi(s_t) < 0$. In order to reduce failure in the early stage of the throwing process, we consider the intrinsic advantage $A_I(s_t, s_{t+1})$ as part of the total advantage function to accelerate the RL training, shown as follows:

$$A_I(s_t, s_{t+1}) = \min(k(V_\phi(s_{t+1}) - V_\phi(s_t)), 0) \quad (9)$$

D. Network Structure

The network architecture for the Lyapunov function is depicted on the right side of Fig. 3. This structure comprises three linear layers, with a Batch normalization layer and an ELU activation function appended to the output of each layer. Following these, there is a penultimate normalization layer [29], succeeded by a linear layer to regulate the output dimension. The network structure for the policy and value function are similar to that of the Lyapunov value network.

To sum up, the total advantage A_{all} consists of 3 parts: standard advantage estimation based on PPO, intrinsic advantage estimation, and stability advantage estimation. The formula for the total advantage is expressed as follows:

$$\begin{aligned} & A_{all}(s_t, a_t, s_{t+1}) \\ = & A^{\pi_\theta}(s_t, a_t) + \beta_1 A_I(s_t, s_{t+1}) + \beta_2 A_L(s_t, a_t) \end{aligned} \quad (10)$$

where $A^{\pi_\theta}(s_t, a_t)$ refers to the time differential advantage function of the PPO algorithm, $A_I(s_t, s_{t+1})$ represents the intrinsic advantage function, and $\beta_1, \beta_2 \in [0, 1]$ are hyperparameters, $A_L(s_t, a_t)$ represents the Lyapunov stability advantage function.

V. EXPERIMENTS

To test the algorithm's ability to generalize, we trained and tested it with objects from the YCB_Pybullet dataset [30]. The results show that our algorithm outperforms baseline methods. Four essential components of our algorithm significantly influence its performance. Additionally, our method showcases the ability to generalize different gestures across three skillful throwing and catching tasks, effectively allowing for the successful manipulation of various objects. Importantly, our approach extends its capabilities to successfully catch previously unseen objects.

A. Comparison Experiments

In this section, we conducted a comparative analysis of our algorithm against mainstream reinforcement learning (RL) algorithms in dexterous hands-throwing and catching tasks. These algorithms include TRPO [31], SAC [32], DDPG [33], and PPO [26]. The results, as illustrated in Fig. 6, are quite telling. In the most challenging task, the Overarm Catch with vertical hands, our algorithm significantly outperforms the baseline methods, achieving success where others fail. For the Abreast Catch task, our approach demonstrates superior learning efficiency in acquiring throwing and catching skills. Furthermore, even in the relatively easier Underarm Catch task, our method exhibits slightly improved learning efficiency and performance compared to the baseline algorithms. These findings highlight the robustness and effectiveness of our algorithm, particularly in demanding scenarios. In the Underarm Catch task, as depicted in Fig. 5, the distance

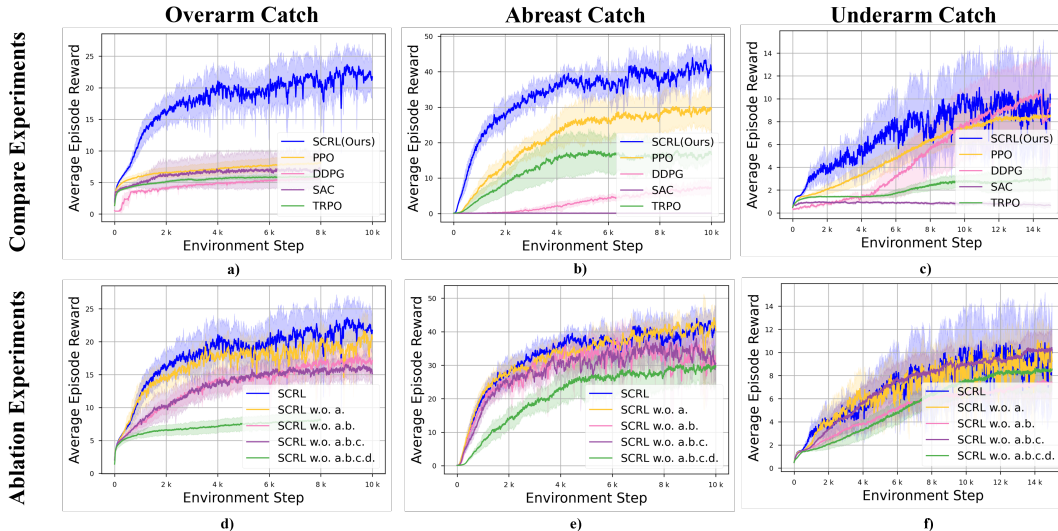


Fig. 6. The average reward curve for training multiple objects for three different tasks. The first row is the comparison experiment reward graph. The second row is an ablation experiment reward graph. (Curves smoothed report the mean 5 times, and shadow areas show the variance.)

TABLE II
SUCCESS RATE OF TASK

Task	Success Rate (%)														
	Train										Test				
	Apple	Banana	Peach	Pen	Strawberry	Pie	Poker	Sugar	Mug	META	Stapler	Scissor	Washer	Bowl	META
Overarm	42.55	53.82	37.20	68.30	28.56	52.71	41.77	48.41	32.20	45.05	23.44	28.40	76.28	82.22	52.58
Abreast	62.11	88.65	46.51	47.12	50.89	99.77	86.34	93.33	76.68	72.37	94.84	27.59	83.94	83.59	72.49
Underarm	71.51	38.80	73.28	31.81	44.62	69.22	42.62	26.34	76.26	52.71	71.33	31.80	57.67	53.24	53.51

between the two hands is the shortest. Consequently, both the PPO and DDPG algorithms manage to solve this task, with DDPG even surpassing our algorithm in terms of average reward. However, it's essential to note that during testing, DDPG and PPO produce less favorable results compared to our method. This discrepancy can be attributed to the fact that their strategies generate unstable grasping gestures, primarily due to the absence of stability requirements in their approaches. Our algorithm, with its focus on stability, excels in ensuring consistent and reliable grasping gestures, even in challenging scenarios. In summary, our algorithm exhibits superior learning efficiency and stability. Video snapshots illustrating the performance in three tasks can be seen in Fig. 5. It's worth emphasizing that we trained the model using nine objects with diverse shapes and mass distributions. This showcases the algorithm's robust generalization ability, as it can effectively adapt to various objects and perform tasks efficiently and reliably.

B. Ablation Experiment

We evaluated the impact of four crucial components in our algorithm. To clarify, here are their definitions: *a.* includes point cloud features from different objects as input. *b.* represents the Lyapunov stability component. *c.* stands for the intrinsic advantage component. *d.* indicates parameter nor-

malization. The effect of these components is shown in the reward curve in Figure 6. Our initial experiment involved the removal of component *a.*, which entails excluding the point cloud features of objects as input. The outcome demonstrated a substantial decrease in the average reward, particularly in the challenging Overarm Catch task. This result underscores the importance of providing the system with prior information about objects to facilitate faster learning in demanding tasks. In our next experiment, we took out components *a.* and *b.* from the algorithm. The reward curve showed a significant effect of the Lyapunov stability component on early learning and how quickly the algorithm converges. This part adds constraints to the learning process, ensuring that the dexterous hands focus on learning stable gestures for catching objects rather than just pursuing the highest rewards. This balance between stability and optimization is essential for the system's performance. For training w.o. *a. b. c.*, it can be seen that the effect of the intrinsic advantage part is slight in Overarm Catch and Abreast Catch. Nevertheless, it offers an important advantage in performance in Underarm Catch. The experiment where we removed all components (*a.*, *b.*, *c.*, and *d.*) underscores the significance of parameter normalization. Notably, this experiment reveals that each element in our algorithm has a slight impact on the Underarm Catch task. This observation aligns with the understanding

that the Underarm Catch task is relatively easier compared to the other tasks.

C. Generalization Ability

For diverse objects: In this section, we delve into a detailed analysis of our method’s generalization capabilities. It demonstrates the ability to execute distinct catching gestures based on the positions of the shadow hands and the objects being caught. Through an analysis of the high-dimensional point cloud features, we categorized most objects into three primary classes. Remarkably, our method exhibits over three different catching gestures for each of these three object types, as depicted in Fig. 7. For objects falling into the regular, flat, and elongated categories, the model generalizes three distinct gestures in each of the three tasks. For instance, when dealing with regular objects like apples, it employs either a half-grip or pinching gesture. Slender objects like bananas are handled with a full-grip position. Flat objects, such as scissors and pokers, are directly held and secured between the fingers and palms. This versatility in catching gestures reflects the algorithm’s remarkable adaptability to different object types and positions.

For different initial postures: Furthermore, given that the initial position of the object is randomly selected in each episode, the object’s motion follows different trajectories due to aerodynamic processes. Impressively, our method excels at solving dynamic catching tasks by generating diverse catching gestures, as illustrated in Fig. 8. For example, when capturing a pen, we observe that the catching hand can adapt to grasp the pen with varying postures. In practical scenarios, catching a randomly rotating pen in mid-air, especially due to its anisotropic shape, would be challenging. However, our results demonstrate the algorithm’s ability to catch randomly rotating objects in the air, even when dealing with challenging objects like pens.

For unseen objects: Moreover, we evaluate our method in the throwing and catching task using various objects, quantifying the success rate as a key metric. We run the testing task in 10k episodes and use the average reward threshold to determine whether it is successful. Since our algorithm pays more attention to the efficiency and stability of generalization ability learning, we no longer use the distance threshold of the final state from the target position as the success rate. The success rate of 13 objects is illustrated in the Table II. Among them, we use 4 unseen objects (test set) to verify the generalization ability of our method. Specifically, we adopt two methods to test objects in the training and test sets: i. Test the success rate of multiple objects in parallel. ii. Test the success rate of each object individually. We take the most challenging task, Overarm Catch as an example. The average success rate for the training set objects reaches 45.05%, and the highest of individual tests on each object was pen up to 68.30%. In addition, the average success rate of the test set object is close to 52.58%. Most importantly, our method obtains a great success rate, 28.40%, in testing scissors, which is a new object. Furthermore, we can observe that our method can provide a robust catching strategy because

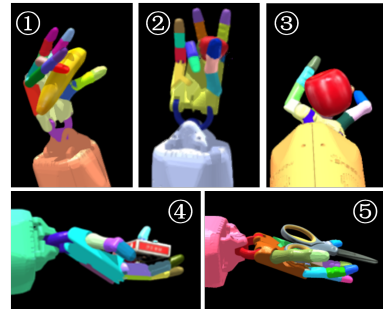


Fig. 7. For the regular, flat, and elongated classes of objects, the model generalizes three different stable catching strategies for unseen objects in each of the three tasks. For regular objects such as apples and strawberries, it uses half-grip (③) or pinching gestures (②). A full-grip gesture (①) is used for slender objects such as bananas. For flat objects such as scissors and pokers, the hand can hold them directly in the palm (④,⑤).

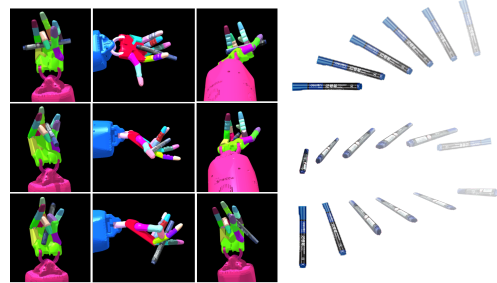


Fig. 8. When the pen is oriented obliquely or tumbling because of the randomness in the initial posture, the shadow hand can generalize various grasping gestures to successfully complete the tasks.

there is not a significant gap between the success rates of training and test objects. For Abreast Catch and Underarm Catch, our method achieves higher success rates in single and multiple objects tests. The results provide strong evidence of our algorithm’s robust generalization capability for catching diverse objects.

D. Implementation detail

In our experiments, we conducted training and testing on a computing platform equipped with an RTX A6000 GPU and an AMD CPU. We utilized nine objects for training across three different task types and tested the algorithm with four objects. Our approach also leveraged meta-learning for multi-object parallel training. The objects used in our experiments were sourced from the YCB dataset [30]. Due to the limitations of the training platform, we employed 256 parallel environments per object for the Overarm Catch and Underarm Catch tasks. For the Abreast Catch task, each object was trained with 64 parallel environments. The policy updates occurred at a frequency of 0.125 times per environment step, and the length of each episode was set to 100 steps.

VI. CONCLUSION

In this paper, we propose the Stability Constraint Reinforcement Learning (SCRL) algorithm to tackle the dexterous shadow-handed throwing and catching task across

diverse objects. We evaluate its performance in Underarm Catch, Abreast Catch and Overarm Catch and find that the SCRL algorithm significantly outperforms baselines. Notably, our learned strategy exhibits remarkable zero-shot transfer capabilities when dealing with previously unseen objects. Even in challenging scenarios where the hand is oriented vertically, our algorithm excels in baselines by a large margin. The incorporation of Lyapunov stability constraints plays a pivotal role in enhancing the stability of the throwing and catching actions. The intrinsic advantage and the regularization of the parameter space mechanism accelerate early learning, leading to a higher performance. Our method autonomously generates different catching gestures for slender, flat, and regular objects. Additionally, it creates diverse catching behaviours when faced with different postures of a type of object. This research opens up promising avenues for the generalization of dynamic manipulation using dexterous hands, showcasing the potential of the RL-based approach to improve object transportation efficiency through throwing and catching.

REFERENCES

- [1] T. Asfour, P. Azad, N. Vahrenkamp, K. Regenstein, A. Bierbaum, K. Welke, J. Schroeder, and R. Dillmann, "Toward humanoid manipulation in human-centred environments," *Robotics and Autonomous Systems*, vol. 56, no. 1, pp. 54–65, 2008.
- [2] U. Kim, D. Jung, H. Jeong, J. Park, H.-M. Jung, J. Cheong, H. R. Choi, H. Do, and C. Park, "Integrated linkage-driven dexterous anthropomorphic robotic hand," *Nature communications*, vol. 12, no. 1, p. 7177, 2021.
- [3] Y. Bai and C. K. Liu, "Dexterous manipulation using both palm and fingers," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1560–1565, IEEE, 2014.
- [4] V. Kumar, E. Todorov, and S. Levine, "Optimal control with learned local models: Application to dexterous manipulation," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 378–383, IEEE, 2016.
- [5] A. M. Okamura, N. Smaby, and M. R. Cutkosky, "An overview of dexterous manipulation," in *Proceedings 2000 ICRA. Millennium Conference. IEEE International Conference on Robotics and Automation. Symposia Proceedings (Cat. No. 00CH37065)*, vol. 1, pp. 255–262, IEEE, 2000.
- [6] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter, "Learning agile and dynamic motor skills for legged robots," *Science Robotics*, vol. 4, no. 26, p. eaau5872, 2019.
- [7] P. Wu, A. Escontrela, D. Hafner, P. Abbeel, and K. Goldberg, "Daydreamer: World models for physical robot learning," in *Conference on Robot Learning*, pp. 2226–2240, PMLR, 2023.
- [8] L. Smith, I. Kostrikov, and S. Levine, "A walk in the park: Learning to walk in 20 minutes with model-free reinforcement learning," *arXiv preprint arXiv:2208.07860*, 2022.
- [9] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning robust perceptive locomotion for quadrupedal robots in the wild," *Science Robotics*, vol. 7, no. 62, p. eabk2822, 2022.
- [10] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine, "Learning complex dexterous manipulation with deep reinforcement learning and demonstrations," *Robotics: Science and Systems XIV*, 2018.
- [11] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, *et al.*, "Learning dexterous in-hand manipulation," *The International Journal of Robotics Research*, vol. 39, no. 1, pp. 3–20, 2020.
- [12] I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, *et al.*, "Solving rubik's cube with a robot hand," *arXiv preprint arXiv:1910.07113*, 2019.
- [13] T. Chen, J. Xu, and P. Agrawal, "A system for general in-hand object re-orientation," in *Conference on Robot Learning*, pp. 297–307, PMLR, 2022.
- [14] R. Jeong, J. T. Springenberg, J. Kay, D. Zheng, Y. Zhou, A. Galashov, N. Heess, and F. Nori, "Learning dexterous manipulation from sub-optimal experts," *arXiv preprint arXiv:2010.08587*, 2020.
- [15] A. Zeng, S. Song, J. Lee, A. Rodriguez, and T. Funkhouser, "Tossing-bot: Learning to throw arbitrary objects with residual physics," *IEEE Transactions on Robotics*, vol. 36, no. 4, pp. 1307–1319, 2020.
- [16] B. Bäuml, T. Wimböck, and G. Hirzinger, "Kinematically optimal catching a flying ball with a hand-arm-system," in *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 2592–2599, IEEE, 2010.
- [17] A. Namiki, Y. Imai, M. Ishikawa, and M. Kaneko, "Development of a high-speed multifingered hand system and its application to catching," in *Proceedings 2003 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2003)(Cat. No. 03CH37453)*, vol. 3, pp. 2666–2671, IEEE, 2003.
- [18] H. J. Charlesworth and G. Montana, "Solving challenging dexterous manipulation tasks with trajectory optimisation and reinforcement learning," in *International Conference on Machine Learning*, pp. 1496–1506, PMLR, 2021.
- [19] G. Williams, A. Aldrich, and E. Theodorou, "Model predictive path integral control using covariance variable importance sampling," *arXiv preprint arXiv:1509.01149*, 2015.
- [20] W. Huang, I. Mordatch, P. Abbeel, and D. Pathak, "Generalization in dexterous manipulation via geometry-aware multi-task learning," *arXiv preprint arXiv:2111.03062*, 2021.
- [21] A. Petrenko, A. Allshire, G. State, A. Handa, and V. Makoviychuk, "Dexpbt: Scaling up dexterous manipulation for hand-arm systems with population based training," *arXiv preprint arXiv:2305.12127*, 2023.
- [22] Y. Chen, T. Wu, S. Wang, X. Feng, J. Jiang, Z. Lu, S. McAleer, H. Dong, S.-C. Zhu, and Y. Yang, "Towards human-level bimanual dexterous manipulation with reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 35, pp. 5150–5163, 2022.
- [23] K. Zakka, L. Smith, N. Gileadi, T. Howell, X. B. Peng, S. Singh, Y. Tassa, P. Florence, A. Zeng, and P. Abbeel, "Robopianist: A benchmark for high-dimensional robot control," *arXiv preprint arXiv:2304.04150*, 2023.
- [24] X. Zhu, R. Zhang, B. He, Z. Zeng, S. Zhang, and P. Gao, "Pointclip v2: Adapting clip for powerful 3d open-world learning," *arXiv preprint arXiv:2211.11682*, 2022.
- [25] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, *et al.*, "Isaac gym: High performance gpu-based physics simulation for robot learning," *arXiv preprint arXiv:2108.10470*, 2021.
- [26] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [27] Y.-C. Chang and S. Gao, "Stabilizing neural control using self-learned almost lyapunov critics," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1803–1809, IEEE, 2021.
- [28] Y. Li, T. Gao, J. Yang, H. Xu, and Y. Wu, "Phasic self-imitative reduction for sparse-reward goal-conditioned reinforcement learning," in *International Conference on Machine Learning*, pp. 12765–12781, PMLR, 2022.
- [29] J. Bjorck, C. P. Gomes, and K. Q. Weinberger, "Is high variance unavoidable in rl? a case study in continuous control," in *International Conference on Learning Representations*, 2021.
- [30] B. Calli, A. Singh, J. Bruce, A. Walsman, K. Konolige, S. Srivastava, P. Abbeel, and A. M. Dollar, "Yale-cmu-berkeley dataset for robotic manipulation research," *The International Journal of Robotics Research*, vol. 36, no. 3, pp. 261–268, 2017.
- [31] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International conference on machine learning*, pp. 1889–1897, PMLR, 2015.
- [32] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International conference on machine learning*, pp. 1861–1870, PMLR, 2018.
- [33] 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," *arXiv preprint arXiv:1509.02971*, 2015.