# Dexterous Non-Prehensile Manipulation for Ungraspable Object via Extrinsic Dexterity

**Yuhan Wang**[1,2]**, Yu Li**[1,2]**, Yaodong Yang**[1,2,†] **and Yuanpei Chen**[1,2,†]
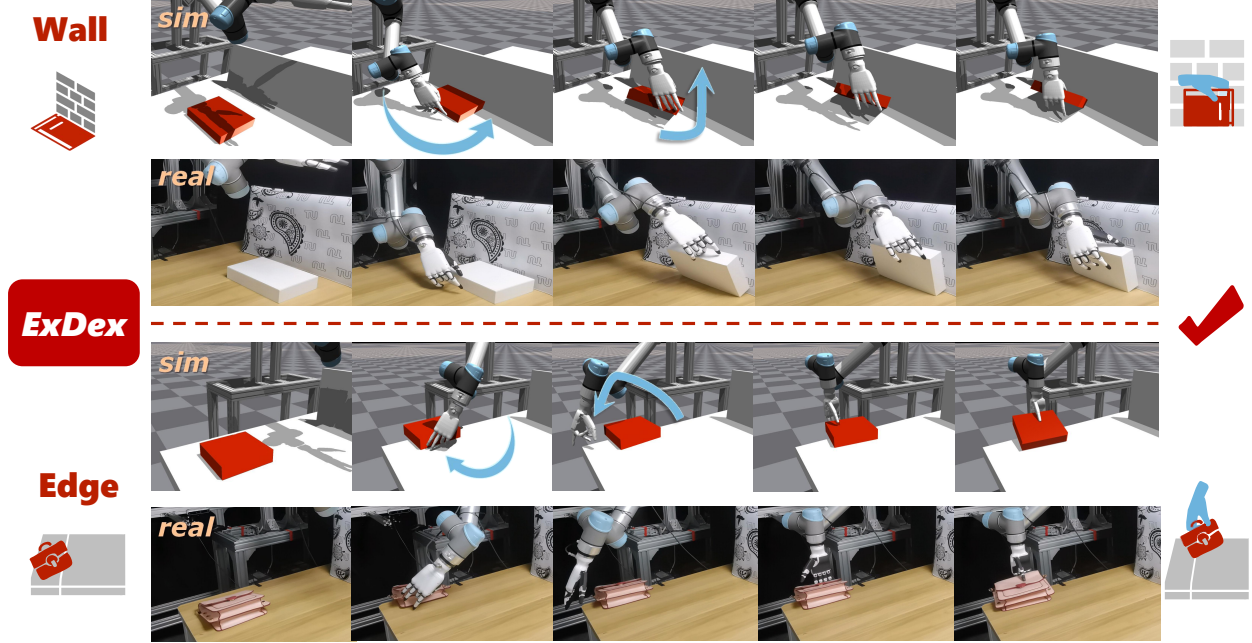[1]Peking University, [2]PKU-PsiBot Joint Lab



Figure 1 | We present **ExDex**, a hierarchical framework that enables non-prehensile manipulation skills for ungraspable objects leveraging external environments. Our approach is demonstrated through two representative tasks: **Wall**, where objects are pushed against walls to facilitate manipulation. And **Edge**, where objects are repositioned to table edges allowing the hand to maneuver into optimal grasping poses. Through reinforcement learning in simulation, we successfully train these manipulation policies and achieve zero-shot transfer to real-world scenarios.

## Abstract

Objects with large base areas become ungraspable when they exceed the end-effector's maximum aperture. Existing approaches address this limitation through extrinsic dexterity, which exploits environmental features for non-prehensile manipulation. While grippers have shown some success in this domain, dexterous hands offer superior flexibility and manipulation capabilities that enable richer environmental interactions, though they present greater control challenges. Here we present ExDex, a dexterous arm-hand system that leverages reinforcement learning to enable non-prehensile manipulation for grasping ungraspable objects. Our system learns two strategic manipulation sequences: relocating objects from table centers to edges for direct grasping, or to walls where extrinsic dexterity enables grasping through environmental interaction. We validate our approach through extensive experiments with dozens of diverse household objects, demonstrating both superior performance and generalization capabilities with novel objects. Furthermore, we successfully transfer the learned policies from simulation to a real-world robot system without additional training, further demonstrating its applicability in real-world scenarios. Project website: https://tangty11.git hub.io/ExDex/.

## 1. Introduction

Humans naturally manipulate objects with their multi-finger dexterous hands through a rich repertoire of strategies. Beyond direct grasping, humans demonstrate remarkable ability to exploit environmental features for manipulation. For instance, when encountering large, flat objects placed in the middle of a table that are challenging to grasp directly, humans intuitively leverage environmental constraints like walls or table edges. They seamlessly combine non-prehensile actions such as pushing, sliding, and pivoting with dexterous manipulation to achieve reliable grasps. This adaptive exploitation of environmental affordances enables humans to handle objects that would oth-

arXiv:2503.23120v1 [cs.RO] 29 Mar 2025

erwise be ungraspable through direct manipulation alone. Such environment-aware manipulation strategies significantly expand the range of objects that can be successfully manipulated, demonstrating the sophisticated interplay between dexterous control and environmental interaction.

Replicating human-like extrinsic dexterity in multi-finger robotic hands remains an unexplored yet crucial problem in robotics. Traditional approaches to dexterous manipulation primarily rely on trajectory optimization with simplified contact models [1, 2]. However, these methods often fail in contact-rich scenarios due to the complexity of modeling dynamic contact interactions and the uncertainty in physical parameters. While imitation learning has demonstrated promising results in direct dexterous manipulation tasks [3, 4, 5], it faces significant limitations when applied to extrinsic dexterity. The collection of high-quality demonstration data through human teleoperation becomes particularly challenging for dynamic contact-rich manipulations, as operators struggle to precisely control multiple fingers while maintaining stable environmental contacts. Recent years have witnessed remarkable progress in applying reinforcement learning to robotic systems [6, 7, 8, 9]. Leveraging large-scale parallel simulation, reinforcement learning enables robots to undergo extensive training in simulation and then deploy the learned policies to the real-world. This approach offers several key advantages for learning extrinsic dexterity skills: First, parallel simulation environments allow for rapid policy exploration without concerns about physical robot wear or safety constraints. Second, we can systematically vary object properties and environmental conditions during training, leading to more robust policies. Third, the framework naturally handles contact-rich scenarios without requiring explicit contact modeling or expert demonstrations. These benefits make reinforcement learning particularly suitable for learning complex manipulation strategies that combine dexterous control with environmental interactions.

Unlike previous works that train end-to-end reinforcement learning policies [10], relocating objects from table centers to edges or walls requires a more sophisticated approach. This task demands a sequence of non-prehensile manipulation skills and strategic planning. While existing research often simplifies the problem by placing objects near external contacts [11, 12], they overlook the complexity of strategic object repositioning and environmental interaction. Even training individual manipulation skills through reinforcement learning presents significant challenges. The high-dimensional action space of multi-finger dexterous hands, combined with the contact-rich nature of environmental interactions, makes each subtask difficult

to learn. Beyond these difficulties, the challenge lies in high-level strategic planning. Optimal object relocation requires consideration of multiple factors: the current object position, robot arm configuration, and kinematic constraints. Simply choosing the nearest environmental contact point is insufficient and often leads to suboptimal or failed manipulations. Instead, the system must evaluate potential target locations while considering the robot's reachability, joint limits, and possible collision-free paths. This comprehensive planning approach significantly improves manipulation success rates compared to naive nearest-neighbor strategies.

Here we present **ExDex** a framework for dexterous manipulation of ungraspable objects using extrinsic dexterity with multi-finger hands, focusing particularly on leveraging table edges and walls. We introduce a hierarchical learning approach combining a high-level planner for identifying optimal environmental contacts with a low-level controller for precise non-prehensile manipulation. The high-level planner generates target poses and transition signals, while the low-level controller executes pushing policies to achieve these poses, followed by object retrieval under external contacts. The experiments in both simulation and real-world settings validate our framework's effectiveness. Our results demonstrate successful generalization to unseen objects and zero-shot transfer to physical systems.

In summary, our main contribution encompasses:

- First exploration of extrinsic dexterity with multi-finger dexterous manipulation in both simulation and real-world scenarios.
- Novel hierarchical framework combining high-level planning and low-level control for occluded grasp tasks.
- Extensive experimental validation demonstrating system effectiveness across simulated and physical environments.

## 2. Related Works

### 2.1. Dexterous Manipulation

Multi-finger dexterous manipulation remains a significant challenge in robotics. Traditional approaches attempt to solve this through trajectory optimization based on dynamic models [1, 2, 13, 14]. However, these models often rely on simplified contact assumptions, failing when confronted with complex, contact-rich tasks. While [12] proposed combining graph search, optimal control, and learning-based objective functions for extrinsic dexterity in pre-grasp operations, their results were limited to simulation, with some predicted trajectories proving impractical for physical robot execution. Recent research has

demonstrated remarkable success with imitation learning [15, 3, 4, 5, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30]. 3D-ViTac [31] achieves precise manipulation using tactile feedback, while DexCap [5] enables complex bimanual tasks through in-the-wild data collection using data gloves. However, imitation learning's reliance on human demonstration data presents substantial costs. Moreover, teleoperation latency makes collecting data for highly dynamic actions particularly challenging (such as flipping objects from wall edges).

In recent years, reinforcement learning (RL) has been widely adopted for dexterous hand manipulation tasks, spanning in-hand object reorientation [10, 32, 33, 34, 35, 36, 6, 7, 8, 9, 37], bimanual manipulation [38, 39, 40], pre-grasping [11, 41], hand long-horizon manipulation [42, 38]. Dynamic Handover [38] demonstrated multi-agent reinforcement learning for dynamic ball-catching between dexterous hands, while Robopianist [43] achieved human-level piano playing through simulation-based training. We develop policies that adapt to dynamic motion control of real robots using RL. To our knowledge, this work represents the first exploration of extrinsic dexterity with dexterous hands demonstrated in both simulation and real-world environments.

## 2.2. Extrinsic Dexterity

External environmental resources such as contacts, gravity, and dynamic motions[44] enable robot hands to grasp and manipulate objects even without suitable contact points. Previous work has demonstrated the utility of environmental interactions, including external contacts for object grasping[11, 45, 41, 12] and reorientation [46], as well as leveraging gravity [47] or dynamic motions [48, 44] to improve grasping postures. However, except for [12], these works primarily employ grippers or underactuated multi-finger hands. We instead utilize a five-fingered dexterous hand, leveraging its greater degrees of freedom and flexibility for enhanced grasping capabilities and improved policy generalization across multiple objects.

Zhou et al. [11] present the closest approach to ours, learning a closed-loop policy through reinforcement learning. However, their work relies on restrictive assumptions, including objects being initially positioned near walls and walls being sufficiently low for grippers to access objects from above. In contrast, our approach accommodates objects anywhere in the workspace. PreAfford [41] learns affordance to guide object-centric interaction point selection for extrinsic dexterity and object repositioning. However, their reliance on motion planning or assumed availability of low-level controllers becomes problematic in complex environments or with challenging objects. Moreover,

most previous work has not explored how to leverage the high degrees of freedom of dexterous hands for extrinsic dexterity. While UniDexFPM [49] investigated dexterous hand pre-grasp manipulation in tabletop environments, their results were limited to simulation without physical robot validation.

We first relocate objects to locations where external conditions can be effectively utilized, such as beside walls or table edges, through non-prehensile manipulation (pushing) before leveraging these conditions for object retrieval. To address the challenge of insertion of fingers between walls and objects after pushing, we learn a strategy that maintains appropriate spacing when positioning objects against walls.

## 3. Task Formulation

In this paper, we address the challenge of grasping ungraspable objects that having large, flat base surfaces using a dexterous multi-finger hand. The task objective is to employ non-prehensile manipulation to reposition objects near environmental features that can assist in successful grasping. We formulate this task as a finite horizon Markov Decision Process (MDP), defined by the 5-tuple $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$. Here, $\mathcal{S}$ and $\mathcal{A}$ represent the state and action spaces respectively. The stochastic dynamics $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ determine the probability of transitioning to state $s'$ given current state $s$ and action $a$. $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ defines the reward function, and $\gamma \in (0, 1)$ is the discount factor. Our objective is to train a policy $\pi$ that maximizes the expected cumulative reward $\mathbb{E}_{\pi}[\sum_{t=0}^{T-1} \gamma^t R]$ in an episode with $T$ time steps.

## 4. Method

In this section, we introduce our system for dexterous non-prehensile manipulation for ungraspable object. The overview of the system is shown in Figure 2. Our framework consists of three parts: the High-level Planner Design (Section 4.1), Low-level Policy Training (Section 4.2) and Joint Training(Section 4.3). The details of our sim-to-real policy transfer are introduced in Section 4.4.

### 4.1. High-level Planner Design

The first step in extrinsic dexterity is relocating objects to environments that can be leveraged for manipulation, such as walls or table edges. Therefore, planning a desired location where external conditions can be effectively utilized is crucial for successful extrinsic dexterity. To achieve this, we train a prediction model $\pi_{pre}$ to predict target positions for object relocation. The model takes environmental point cloud data $p$ as input and outputs three-dimensional target coordinates $P_t = (P_x, P_y, P_z)_t$ and a signal $s$ to pick a low-level policy automatically for subsequent applying. The pre-
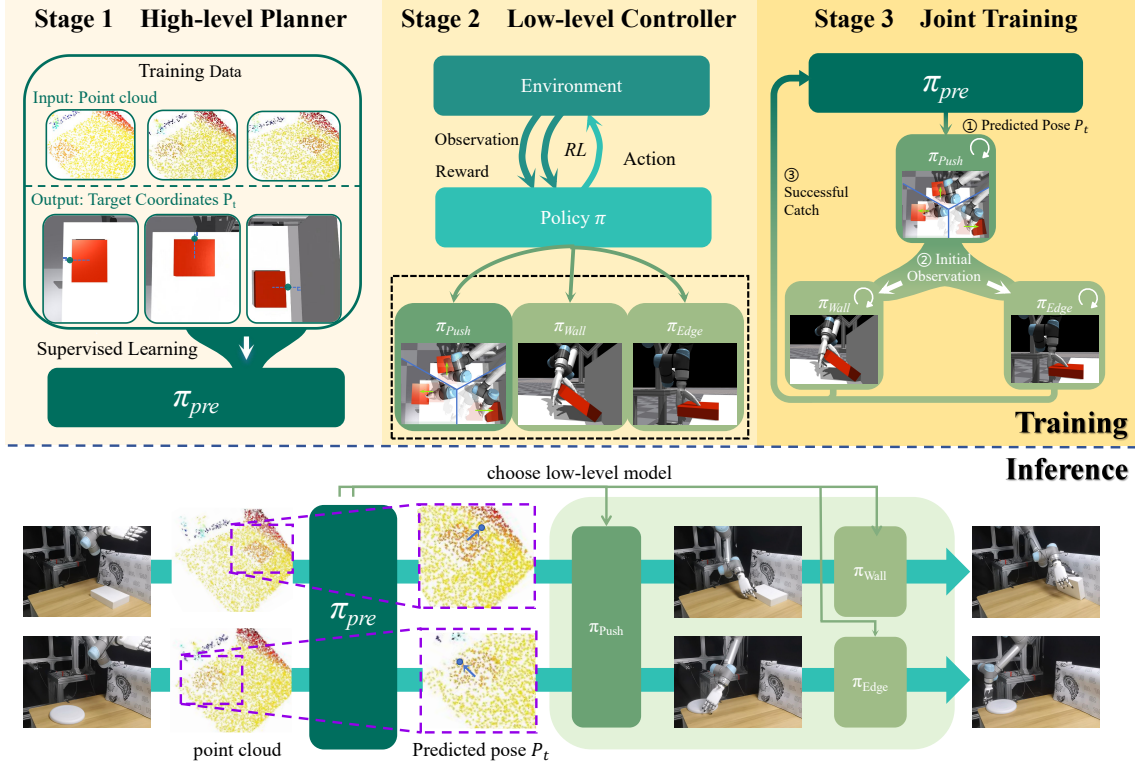
Figure 2 | **Illustration of the ExDex System Design.** (A) Training: Our system is trained in three stages. In Stage 1, we train a prediction model $\pi_{pre}$ through supervised learning that takes point cloud input and predicts the optimal target position $P_t$ for object repositioning. Stage 2 focuses on training three low-level skills via reinforcement learning: a pushing policy $\pi_{push}$ that repositions objects to target locations, and two policies $\pi_{wall}, \pi_{edge}$ that enable grasping of ungraspable objects from walls and table edges via extrinsic dexterity. In Stage 3, we jointly finetune these policies to ensure better transitions between consecutive skills. (B) Inference: During inference, our system first use the $\pi_{pre}$ to process the environmental point cloud to determine whether to execute the $\pi_{wall}$ or $\pi_{edge}$, while simultaneously predicting the corresponding target position $P_t$. The pushing policy $\pi_{push}$ then moves the object to this target position, followed by the selected extrinsic dexterity policy ($\pi_{wall}$ or $\pi_{edge}$) to complete the grasp.

dicted $P_t$ serves as the target position for the $\pi_{push}$ policy in the low-level controller.

## 4.2. Low-level Policy Training

Based on the predicted target pose $P_t$ from the high-level planner, we train three specialized policies using model-free reinforcement learning: (1) A $\pi_{push}$ policy that pushes objects to the target position $P_t$; (2) A $\pi_{wall}$ policy for grasping objects near walls starting from $P_t$; and (3) A $\pi_{edge}$ policy for retrieving objects from table edges at $P_t$. The following subsections detail the observation and action space, reward design, and training methodology.

**Observation Space.** The observation space

$$\mathcal{S} = \{q_t, \{F_t^{f,i}\}_{i=1}^5, p_t^{\text{obj}}, v_t^{\text{obj}}, P_t, c_p\} \qquad (1)$$

consists of several components: robot state (proprioceptive arm and hand joint positions $q_t \in \mathbb{R}^{18}$, five fingertip poses $\{F_t^{f,i}\}_{i=1}^5 \in \mathbb{R}^{15}$), object pose $p_t^{\text{obj}} \in SE(3)$

and velocity $v_t^{\text{obj}} \in \mathbb{R}^6$, target information (predicted pose $P_t$ from high-level planner), and contact information (hand-designed contact position $c_p \in \mathbb{R}^3$ that maintains a fixed relative position to the object center).

**Action Space.** The action space $a_t = \{a_t^{\text{arm}}, a_t^{\text{hand}}\}$ consists of two components: hand joint positions $a_t^{\text{hand}} \in \mathbb{R}^6$ and relative arm joint positions $a_t^{\text{arm}} \in \mathbb{R}^6$. For the hand, the policy directly outputs absolute joint angles $a_t^{\text{hand}}$ as target positions. For the arm, the policy generates relative position changes $a_t^{\text{arm}}$, which are added to the current joint angles to obtain target positions. The PD controller then converts these target positions into joint torques for both the arm and hand.

**Reward Design.** To reduce the complexity of reward shaping, we unify our reward function into three components with a staged reward mechanism. The next stage reward is only calculated when specific condi-

tions are met. Specifically:

$$r = r_{\text{motion}} + r_{\text{pregrasp}} \cdot P(a) + r_{\text{grasp}} \cdot P(b) \qquad (2)$$

where $P(\cdot)$ represents condition probabilities. In the following, we describe each reward component in detail. All reward terms share the same goal of minimizing distances between their arguments, thus we denote these proximity-based functions as $T(\cdot, \cdot)$, which output larger values as their arguments become closer. The specific implementation of $P(\cdot)$, $T(\cdot, \cdot)$ and hyperparameter can be found in Appendix A.1.

### 4.2.1. Motion reward $r_{motion}$

The motion reward is designed to guide either object movement to a target pose or fingertip positioning for manipulation. For $\pi_{\text{push}}$ and $\pi_{\text{wall}}$, it encourages the object to reach specific target poses: $r_{\text{motion}} = T(P_t^{obj}, P_t^{target})$. In $\pi_{\text{push}}$, $P_t^{target}$ is set to the pose $P_t$ predicted by the high-level planner, while in $\pi_{\text{wall}}$, $P_t^{target}$ is a pre-defined pose above the object to facilitate extrinsic dexterity. For $\pi_{\text{edge}}$, the reward guides fingertip positioning: $r_{\text{motion}} = T(\{F_t^{f,i}\}_{i=1}^5, P_t^{target})$, encouraging the thumb to position above the object while placing the other four fingers beneath the object.

### 4.2.2. Pre-grasp reward $r_{pregrasp}$

The pre-grasp reward encourages the hand to achieve an advantageous pre-grasp pose after object repositioning: $r_{\text{pregrasp}} = T(F_t^{f,3}, c_p)$, where $F_t^{f,3}$ is the position of the middle fingertip, and $c_p$ is a relative fixed point on the object.

### 4.2.3. Grasp reward $r_{grasp}$

Once in pre-grasp position, the grasp reward promotes stable grasping by optimizing the collective positions of all five fingertips relative to the object: $r_{\text{grasp}} = T(P_t^m, p_t^{obj})$, where $P_t^m = \frac{F_t^{f,1} + F_t^{f,2}}{3}$ represents the midpoint between thumb ($F_t^{f,1}$) and middle fingertip ($F_t^{f,3}$) positions.

**Training Method** We employ PPO [50] to train our low-level control policies, leveraging its stability and sample efficiency. The training process is accelerated through parallel simulations in IsaacGym, enabling simultaneous training across 4096 environments. To improve policy robustness, we incorporate comprehensive domain randomization techniques, including variations in robot and object properties. Furthermore, we adopt a curriculum learning strategy to enhance training efficiency. The training begins with some similar objects (Figure 3 (a-pretrain)) at a fixed initial position. As the success rate improves, we gradually increase task complexity by introducing objects with larger difference in size(Figure 3 (a-finetune)) and randomizing their initial positions. This progressive

learning approach helps the policies develop robust manipulation skills while maintaining stable training dynamics.

### 4.3. Joint Finetuning

Sequentially executing trained policies often leads to poor performance in long-horizon manipulation tasks. This is primarily because the terminal state of previous policy may not align well with the initial state distribution of the subsequent policy, creating a state distribution mismatch that affects policy execution. This challenge, known as the skill chaining problem [42, 51, 52], requires special consideration in policy training. To address this issue, we jointly finetune our policies with the following order. Firstly, to enhance the robustness of $\pi_{\text{push}}$ against potential biases in the high-level planner's predictions, we introduce Gaussian noise with a specified standard deviation to the predicted target pose. Then, to better align the transition states between sequential non-prehensile skills, we replace the initial states of both $\pi_{\text{wall}}$ and $\pi_{\text{edge}}$ with terminal states obtained from $\pi_{\text{push}}$ rollouts. Finally, to improve the capability of the high level planner $\pi_{\text{pre}}$ to predict the target pose with a higher likelihood of success, we refine it with the target poses from successful rollouts for each task. This approach ensures smoother transitions between different manipulation phases, enhancing the overall task performance.

### 4.4. Sim-to-Real Transfer

When applying the RL policy to the realworld, some environment states can not be estimated accurately like joint and object velocity. Therefore, we follow the teacher-student distilling framework [32, 10] to zero-shot transfer our simulation policy into the real dexterous arm-hand system. Specifically, we rollout our policies in simulation sequentially to collect the whole teacher demonstration trajectories. For distillation from demonstration, we employ a transformer-based imitation learning network to predict the target arm-hand joint angles. To mitigate the observation gap between simulation and real-world, the distilled student policy only takes the low dimension state observation including proprioception and object 6d pose as input. To obtain object 6d pose, we use Segment Anything model[53] to get the initial mask of the object, followed by FoundationPose[54] for pose estimation and tracking. We also build a digital twin framework for sim-to-real transfer. More details about the teacher-student distillation and digital twin can be found in Appendix A.2.

Figure 3 | **Overview of the environment setups.** (a) Object sets used in simulation. Policies are firstly trained on the pretrain set, and then finetuned on the finetune set, and tested for zero-shot generalization on the unseen set. (b) Real-world test objects (top: wall-task objects, bottom: edge-task objects). (c) Workspace of the real-world, We use an Inspired Hand mounted on a UR5e robot, equipped with a RealSense D455 camera.

## 5. Experiment

In this section, we comprehensively evaluate the performance of our proposed framework in simulation and real world to address the following questions:

(1) Can our high-level planner generate optimal object relocating strategy given different external environments? (Section 5.2)

(2) Is the dexterous hand motion learned by our low-level polices necessary for our tasks? (Section 5.3)

(3) Is our reward design suitable for the non-prehensile manipulation skill training? (Section 5.3)

(4) Can our joint finetuning strategy improve the gen-

eralizability and robustness of our framework? (Section 5.3)

(5) Can our framework learned in simulation be applied to a real-world dexterous arm-hand robot system? (Section 5.4)

First, we introduce the main setup of our experiments including our dataset, evaluation metric and several baselines for comparison with our method. Then we evaluate the effectiveness of our framework separately from high-level and low-level parts through quantitative and qualitative results. Finally, we provide details of how we conduct real-world experiments and the performance of our method. Our simulation and real-world settings are shown in Figure 3. All simulation results in tables are evaluated in 3 different seeds, and the real-world results are evaluated in 10 trials for each object.

### 5.1. Setup

**Dataset** In simulation environments, we only use box with various physics property ratios as the training asset. We evaluate the generalizability of our framework on other 15 objects with diverse geometries. In real-world scenario, we select 10 objects with different sizes and physics properties for sim-to-real evaluation. The objects used in simulation and real-world are shown in Figure 3.

**Evaluation and Metric** We evaluate the performance of our framework in a scene as shown in Figure 3(c) containing three external contacts—Wall, FrontEdge and LeftEdge—each representing a separate task. For all tasks, objects are initially placed randomly in the center of the table. We train individual policies for each task as follow. For the Wall task, the policy $\pi_{wall}$ is trained to first push the object to a wall-adjacent position before utilizing the wall to assist in grasping. For the FrontEdge and LeftEdge task, we follow the training paradigm of $\pi_{edge}$, where the object is pushed to the respective table edge prior to maneuvering the hand to an appropriate position for grasping. We introduce the following metrics for evaluating the performance: (1) *Target Transition Error* (TTE) is the Euclidean distance in centimeters between the target position predicted by the high-level planner and ground truth. (2) *Success Rate* (SR) is the percentage of successful grasping after a series of non-prehensile manipulation. We define success if the object is grasped steadily above a height threshold (10 cm). These metrics evaluate different aspects of our framework: TTE assesses the accuracy of the high-level planner, and SR evaluates the overall task completion success.

**Baseline** We compare our methods with the following baselines and ablations. (1) *Random Target*. This
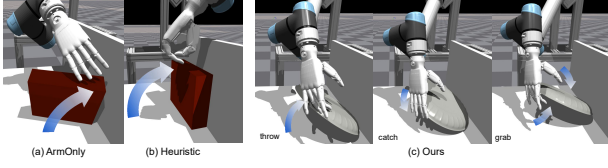
Figure 4 | Comparison of *Arm-Only*, *Heuristic*, and *Ours* for Wall task. (a) *Arm-Only*. (b) *Heuristic*. (c) *Ours*.

baseline maintains our low-level policies but replaces the high-level planner with random target position selection. The random positions are constrained to theoretically feasible regions (e.g., ensuring partial object overhang for the FrontEdge and LeftEdge tasks, or positions near walls for the Wall task) to maintain basic task feasibility. (2) *Arm-Only*. Our low-level policies trained with arm control only. The hand joint angle is fixed. (3) *Heuristic*. Using predefined action primitives as low-level policies. (4) *Ours w/o MR*. Our low-level policies trained without motion reward. (5) *Ours w/o ST*. Our low-level policies trained without stage reward mechanism. (6) *Ours w/o JH*. Our low-level policies without joint finetuning for the high-level planner. (7) *Ours w/o JL*. Our low-level policies without joint finetuning for the low-level policies.

## 5.2. High-level planner

To evaluate the generalizability of the high-level planner for different external contacts, we compare the relocating strategy of our high-level planner with *Random Target*. The quantitative results in Table 1 shows that our high-level planner generates superior relocating positions with lower target transition error (TTE) across all tasks and objects, which facilitates the continual non-prehensile manipulation skills with higher success rate (SR). These results demonstrate that our high-level planner can accurately predict target positions suitable for all the tasks. (add more analysis for different task)Moreover, the improved SR indicates that after joint finetuning, our prediction model tends to select target positions where subsequent policies have higher success rates, suggesting effective integration between the high-level planner and low-level controller.

## 5.3. Low-level controller

**Dexterous hand motion** We compare our method with *Arm-Only* and *Heuristic* to validate the importance of dexterous hand motion for non-prehensile manipulation. As evidenced in Table 2, our approach demonstrates consistent superiority over both baseline methods across all task configurations. The most notable performance gap emerges in the Wall task, where *Arm-Only* and *Heuristic* exhibit fundamental limita-

tions due to their constrained manipulation strategies. *Arm-Only* can only pivot the object upright like Fig 4(a), causing failure in unseen non box-shaped objects. *Heuristic* follows a circular arc trajectory centered on the midpoint of the contact line between the corner and the object. However, it can only rotate the object to squeeze up against the wall, which prevents stable grasping (Fig 4(b)). Our method overcomes these limitations through the learned dexterous hand motion. Specifically, our RL policy $\pi_{\text{wall}}$ leverages finger motions to lift and dynamically catch the object mid-air (Fig. 4(c)), demonstrating superior dexterity. This advantage extends to FrontEdge and LeftEdge task, where our approach maintains robust performance across both seen and unseen objects.

**Reward Design** To investigate the importance of our reward design in low-level policy learning, we conduct ablation studies on two key components: (1) the motion reward (*Ours w/o MR*), which guide the object toward the target pose, and (2) the stage reward mechanism (*Ours w/o ST*), which dynamically adjust different reward components during training. As shown in Table 2, removing the motion reward (*Ours w/o MR*) leads to near zero success rates across all tasks, demonstrating that precise motion reward guidance is essential for low-level policy training. Similarly, ablating the stage reward mechanism (*Ours w/o ST*) causes a drastic performance drop, which confirms that dynamic reward adjustment is critical for smooth transitions between task stages.

**Joint Finetuning** To assess the effectiveness of our joint finetuning approach in enhancing the framework's generalization capability and robustness, we conducted systematic ablation studies examining both the high-level planner (*Ours w/o JH*) and low-level controller (*Ours w/o JL*) components. The experimental results Table 2 reveal that removing joint finetuning for the low-level controller results in substantial performance degradation, particularly on unseen object. Specifically, the success rate drops approaching 40% in these cases, clearly demonstrating that our joint finetuning approach effectively bridge the gap of the state mismatch of chaining low-level policies. While the baseline configuration without high-level planner finetuning (*Ours w/o JH*) maintains reasonable performance across all tasks, our analysis shows that incorporating target poses from successful demonstrations to refine the high-level planner yields consistent slight performance improvements. This suggests that both components of our joint finetuning strategy contribute meaningfully to the framework's overall effectiveness.

## 5.4. Real-world Experiment

**Hardware Setup** We set up the identical scenario in the real world as in the simulation, one multi-finger

Table 1 | Quantitative comparison of the high-level planner.

| Method | Wall | | | | Edge | | | | Left Edge | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Seen | | Unseen | | Seen | | Unseen | | Seen | | Unseen | |
| | TTE | SR | TTE | SR | TTE | SR | TTE | SR | TTE | SR | TTE | SR |
| Random Target | $45.19_{\pm1.94}$ | $66.21_{\pm0.96}$ | $45.15_{\pm1.90}$ | $54.50_{\pm1.47}$ | $31.63_{\pm2.77}$ | $88.23_{\pm1.04}$ | $31.75_{\pm2.87}$ | $60.11_{\pm2.21}$ | $31.23_{\pm1.23}$ | $74.28_{\pm0.12}$ | $31.32_{\pm1.24}$ | $\mathbf{57.22}_{\pm2.70}$ |
| **Ours** | $\mathbf{0.41}_{\pm0.00}$ | $\mathbf{83.25}_{\pm0.34}$ | $\mathbf{0.41}_{\pm0.01}$ | $\mathbf{54.94}_{\pm1.57}$ | $\mathbf{3.16}_{\pm0.03}$ | $\mathbf{89.43}_{\pm0.68}$ | $\mathbf{2.85}_{\pm0.06}$ | $\mathbf{68.00}_{\pm3.55}$ | $\mathbf{2.58}_{\pm0.01}$ | $\mathbf{76.75}_{\pm2.41}$ | $\mathbf{2.97}_{\pm0.07}$ | $54.39_{\pm2.27}$ |

Table 2 | Quantitative comparison of the low-level controller.

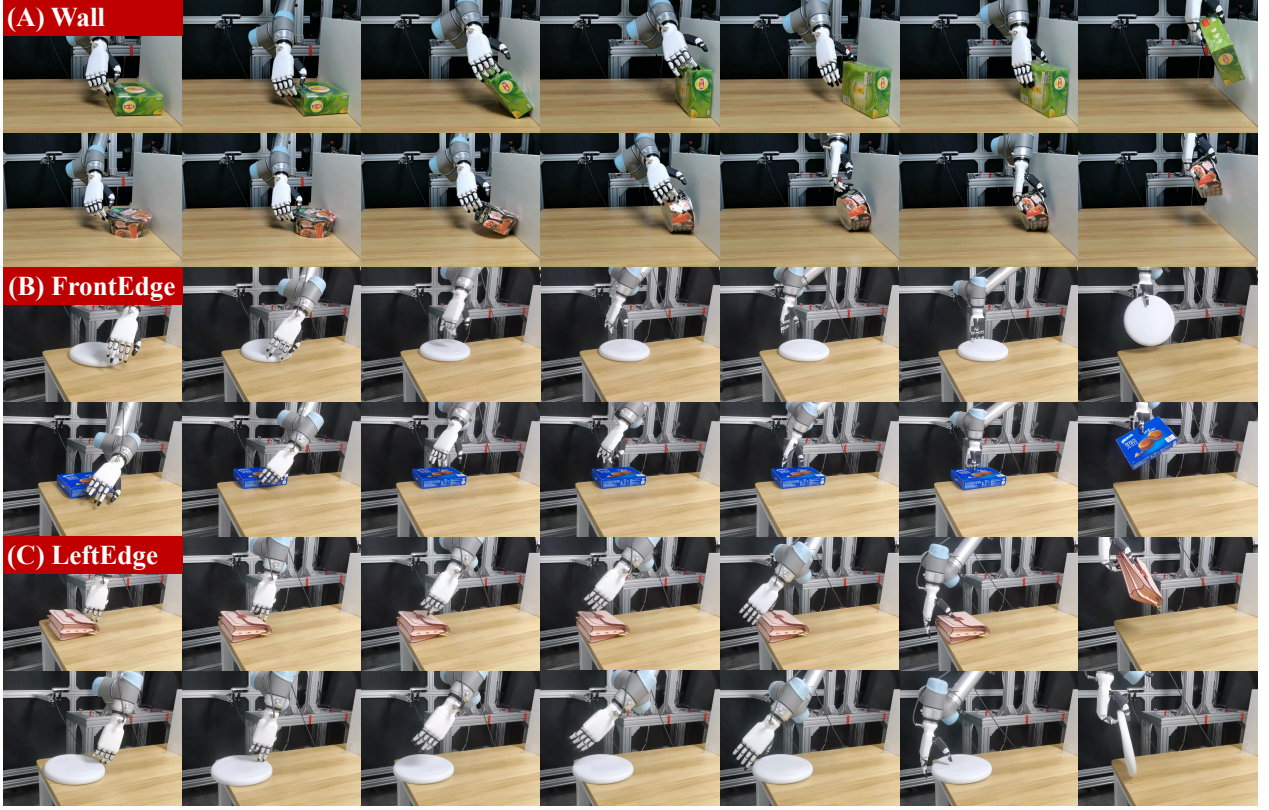| Method | Wall | | FrontEdge | | LeftEdge | |
| --- | --- | --- | --- | --- | --- | --- |
| | Seen | Unseen | Seen | Unseen | Seen | Unseen |
| ArmOnly. | $16.27_{\pm0.18}$ | $2.50_{\pm1.22}$ | $29.95_{\pm0.25}$ | $26.39_{\pm2.74}$ | $14.28_{\pm2.39}$ | $19.61_{\pm0.70}$ |
| Heuristic. | $8.03_{\pm0.30}$ | $2.72_{\pm0.61}$ | $73.70_{\pm0.61}$ | $56.61_{\pm1.86}$ | $63.91_{\pm1.22}$ | $48.61_{\pm0.98}$ |
| w/o MR. | $0.00_{\pm0.00}$ | $0.17_{\pm0.14}$ | $0.00_{\pm0.00}$ | $1.00_{\pm0.14}$ | $0.00_{\pm0.00}$ | $0.01_{\pm0.02}$ |
| w/o ST. | $0.01_{\pm0.02}$ | $0.22_{\pm0.21}$ | $2.28_{\pm0.24}$ | $5.44_{\pm0.42}$ | $0.56_{\pm0.31}$ | $1.36_{\pm0.14}$ |
| w/o JL. | $49.47_{\pm3.67}$ | $19.39_{\pm2.91}$ | $74.56_{\pm3.31}$ | $32.61_{\pm2.21}$ | $48.18_{\pm0.47}$ | $12.94_{\pm0.55}$ |
| w/o JH. | $\mathbf{83.88}_{\pm0.33}$ | $54.28_{\pm0.97}$ | $82.02_{\pm0.92}$ | $63.22_{\pm3.80}$ | $73.33_{\pm2.70}$ | $53.44_{\pm3.91}$ |
| **Ours** | $83.25_{\pm0.34}$ | $\mathbf{54.94}_{\pm1.57}$ | $\mathbf{89.43}_{\pm0.68}$ | $\mathbf{68.00}_{\pm3.55}$ | $\mathbf{76.75}_{\pm2.41}$ | $\mathbf{54.39}_{\pm2.27}$ |



Figure 5 | Real-world experiment demonstrations. The snapshots show successful executions of our framework on various objects. (a) Wall tasks. (b) FrontEdge tasks. (c) LeftEdge tasks.

Inspire Hand mounted on a UR5e robot for our experiments, as shown in Figure 3. To obtain the real-time visual observation for object pose estimation and environment state extraction, we use a single RealSense D455 depth camera.

**Evaluation and Metric** We evaluate our framework in real-world scenarios following the same protocols as our simulation experiments. For all tasks, objects are initially placed randomly in the center of the table. The robot then needs to execute the complete manipulation sequence: predict the target position, pushing

Table 3 | Results for Real-world Experiments using teacher-student distillation

| Size ($cm^3$) | 16.5x15.1x6.2 | 17.3x17.3x7.5 | 23x16.2x5 | 20.7x16.5x7 | 19x14x4 |
|---|---|---|---|---|---|
| Wall | 10/10 | 10/10 | 9/10 | 8/10 | 8/10 |
| Size ($cm^3$) | 21x13x4.4 | 22x16x4.4 | 24.7x23.8x4 | 20.7x16.5x7 | 23.5x23.5x2.3 |
| FrontEdge | 7/10 | 9/10 | 9/10 | 6/10 | 8/10 |
| Size ($cm^3$) | 21x13x4.4 | 22x16x4.4 | 24.7x23.8x4 | 20.7x16.5x7 | 23.5x23.5x2.3 |
| LeftEdge | 5/10 | 7/10 | 8/10 | 5/10 | 7/10 |

the object to the target position near environmental contacts (wall, front or left edge) and subsequently performing the corresponding extrinsic dexterity skill to grasp it. We use the success rate (SR) as our evaluation metric, following the same criteria defined in Section 5.1, where a grasp is considered successful if the object is lifted steadily above a specified height threshold.

**Object set** To thoroughly evaluate the sim-to-real transfer capability of our framework, we conduct experiments on a diverse set of real-world objects that vary significantly in their physical properties. Our test objects, illustrated in Figure 3, include items with different geometries, sizes, and masses. Moreover, we deliberately include several deformable objects, which present additional challenges for non-prehensile manipulation due to their changing dynamics during interaction. This diverse object set allows us to assess the robustness of our policies when dealing with physical properties that are challenging to model accurately in simulation.

**Sim-to-real performance** As shown in Table 5, our policies achieves robust performance on real-world objects, with success rates exceeding 80% for most tasks, demonstrating effective sim-to-real transfer capability. More importantly, our framework maintains high success rates even when handling objects that differ significantly from the simulation training set in terms of size and physical properties. This robust performance extends to challenging scenarios involving deformable objects, whose dynamics are particularly difficult to model accurately in simulation. These results highlight the strong generalization capability and robustness of our trained policies, successfully bridging the reality gap in complex manipulation tasks. The detailed visualization of our real-world experiments are shown in Figure 5.

## 6. Limitation

There are several limitations of our work. (1) *Limited operation space*. Our current implementation relies on a fixed-base robot arm, which constrains the operational workspace to a corner region of the table. This limitation restricts the position generalizability of our policies, as the robot cannot access objects placed beyond its reachable workspace. A potential solution would be to integrate our framework with a mobile manipulator, which would enable extrinsic dexterity capabilities across larger working areas. (2) *Clutter scene generalizability*. All of our experiments are conducted on a clean table which only contains wall or target objects. However, a common table in daily life is usually filled with cluttered objects, which not only introduces additional obstacles for manipulation but also provides potential new external contacts for extrinsic dexterity. Future work could focus on enhancing both our prediction model and pushing policy to enable robust object repositioning in cluttered environments, effectively identifying and utilizing suitable positions for extrinsic dexterity among obstacles.

## 7. Conclusion

In this work, we investigate the challenging problem of manipulating ungraspable objects using extrinsic dexterity with a multi-finger hand. Inspired by human's ability to leverage environmental features like walls and edges, we present a hierarchical framework that combines strategic planning with dexterous manipulation skills. Our framework features a high-level planner that intelligently selects optimal external contacts and predicts target positions, coupled with a low-level controller that executes precise non-prehensile manipulation skills. Through extensive experiments in simulation, we demonstrate our framework's superior performance across different external contacts and various objects. The results show that our approach successfully addresses the key challenges in extrinsic dexterity by three factors: strategic object repositioning, dynamic contact interactions, and precise manipulation control. Moreover, the successful transfer of our policies from simulation to a real-world robot system validates the practical applicability of our method, bridging the gap between simulation and reality in contact-rich manipulation tasks.

## References

[1] Sirui Chen, Jeannette Bohg, and C Karen Liu. Springgrasp: An optimization pipeline for robust and compliant dexterous pre-grasp synthesis. *arXiv preprint arXiv:2404.13532*, 2024.

[2] Igor Mordatch, Zoran Popović, and Emanuel Todorov. Contact-invariant optimization for hand manipulation. In *Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer animation*, pages 137–144, 2012.

[3] Zoey Qiuyu Chen, Karl Van Wyk, Yu-Wei Chao, Wei Yang, Arsalan Mousavian, Abhishek Gupta, and Dieter Fox. Learning robust real-world dex-

terous grasping policies via implicit shape augmentation. *arXiv preprint arXiv:2210.13638*, 2022.

[4] Kenneth Shaw, Shikhar Bahl, and Deepak Pathak. Videodex: Learning dexterity from internet videos. In *Conference on Robot Learning*, pages 654–665. PMLR, 2023.

[5] Chen Wang, Haochen Shi, Weizhuo Wang, Ruohan Zhang, Li Fei-Fei, and C Karen Liu. Dexcap: Scalable and portable mocap data collection system for dexterous manipulation. *arXiv preprint arXiv:2403.07788*, 2024.

[6] OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider, Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba, and Lei Zhang. Solving rubik's cube with a robot hand. *CoRR*, abs/1910.07113, 2019.

[7] Max Yang, Chenghua Lu, Alex Church, Yijiong Lin, Chris Ford, Haoran Li, Efi Psomopoulou, David AW Barton, and Nathan F Lepora. Anyrotate: Gravity-invariant in-hand object rotation with sim-to-real touch. *arXiv preprint arXiv:2405.07391*, 2024.

[8] Johannes Pitz, Lennart Röstel, Leon Sievers, and Berthold Bäuml. Dextrous tactile in-hand manipulation using a modular reinforcement learning architecture. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1852–1858. IEEE, 2023.

[9] Ankur Handa, Arthur Allshire, Viktor Makoviychuk, Aleksei Petrenko, Ritvik Singh, Jingzhou Liu, Denys Makoviichuk, Karl Van Wyk, Alexander Zhurkevich, Balakumar Sundaralingam, et al. Dextreme: Transfer of agile in-hand manipulation from simulation to reality. *arXiv preprint arXiv:2210.13702*, 2022.

[10] Tao Chen, Jie Xu, and Pulkit Agrawal. A system for general in-hand object re-orientation. *Conference on Robot Learning*, 2021.

[11] Wenxuan Zhou and David Held. Learning to grasp the ungraspable with emergent extrinsic dexterity. In *Conference on Robot Learning*, pages 150–160. PMLR, 2023.

[12] Sirui Chen, Albert Wu, and C Karen Liu. Synthesizing dexterous nonprehensile pregrasp for ungraspable objects. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–10, 2023.

[13] Yunfei Bai and C Karen Liu. Dexterous manipulation using both palm and fingers. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1560–1565. IEEE, 2014.

[14] Vikash Kumar, Yuval Tassa, Tom Erez, and Emanuel Todorov. Real-time behaviour synthesis for dynamic hand-manipulation. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6808–6815. IEEE, 2014.

[15] Priyanka Mandikal and Kristen Grauman. Learning dexterous grasping with object-centric visual affordances. In *2021 IEEE international conference on robotics and automation (ICRA)*, pages 6169–6176. IEEE, 2021.

[16] Zoey Qiuyu Chen, Karl Van Wyk, Yu-Wei Chao, Wei Yang, Arsalan Mousavian, Abhishek Gupta, and Dieter Fox. Dextransfer: Real world multi-fingered dexterous grasping with minimal human demonstrations. *arXiv preprint arXiv:2209.14284*, 2022.

[17] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.

[18] Ilija Radosavovic, Xiaolong Wang, Lerrel Pinto, and Jitendra Malik. State-only imitation learning for dexterous manipulation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7865–7871. IEEE, 2021.

[19] Sridhar Pandian Arunachalam, Irmak Güzey, Soumith Chintala, and Lerrel Pinto. Holo-dex: Teaching dexterity with immersive mixed reality. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5962–5969. IEEE, 2023.

[20] Irmak Guzey, Ben Evans, Soumith Chintala, and Lerrel Pinto. Dexterity from touch: Self-supervised pre-training of tactile representations with robotic play. *arXiv preprint arXiv:2303.12076*, 2023.

[21] Ankur Handa, Karl Van Wyk, Wei Yang, Jacky Liang, Yu-Wei Chao, Qian Wan, Stan Birchfield, Nathan Ratliff, and Dieter Fox. Dexpilot: Vision-based teleoperation of dexterous robotic hand-arm system. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9164–9170. IEEE, 2020.

[22] Aravind Sivakumar, Kenneth Shaw, and Deepak Pathak. Robotic telekinesis: Learning a robotic hand imitator by watching humans on youtube. *arXiv preprint arXiv:2202.10448*, 2022.

[23] Yuzhe Qin, Wei Yang, Binghao Huang, Karl Van Wyk, Hao Su, Xiaolong Wang, Yu-Wei Chao, and Dietor Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation system. *arXiv preprint arXiv:2307.04577*, 2023.

[24] Yuzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. In *European Conference on Computer Vision*, pages 570–587. Springer, 2022.

[25] Zichen Jeff Cui, Yibin Wang, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. From play to policy: Conditional behavior generation from uncurated robot data. *arXiv e-prints*, pages arXiv–2210, 2022.

[26] Siddhant Haldar, Jyothish Pari, Anant Rai, and Lerrel Pinto. Teach a robot to fish: Versatile imitation from one minute of demonstrations. *arXiv preprint arXiv:2303.01497*, 2023.

[27] Yuzhe Qin, Hao Su, and Xiaolong Wang. From one hand to multiple hands: Imitation learning for dexterous manipulation from single-camera teleoperation. *IEEE Robotics and Automation Letters*, 7(4):10873–10881, 2022.

[28] Sridhar Pandian Arunachalam, Sneha Silwal, Ben Evans, and Lerrel Pinto. Dexterous imitation made easy: A learning-based framework for efficient dexterous manipulation. *arXiv preprint arXiv:2203.13251*, 2022.

[29] Irmak Guzey, Yinlong Dai, Ben Evans, Soumith Chintala, and Lerrel Pinto. See to touch: Learning tactile dexterity through visual incentives. *arXiv preprint arXiv:2309.12300*, 2023.

[30] Toru Lin, Yu Zhang, Qiyang Li, Haozhi Qi, Brent Yi, Sergey Levine, and Jitendra Malik. Learning visuotactile skills with two multifingered hands. *arXiv preprint arXiv:2404.16823*, 2024.

[31] Binghao Huang, Yixuan Wang, Xinyi Yang, Yiyue Luo, and Yunzhu Li. 3d vitac:learning fine-grained manipulation with visuo-tactile sensing. In *Proceedings of Robotics: Conference on Robot Learning(CoRL)*, 2024.

[32] Tao Chen, Megha Tippur, Siyang Wu, Vikash Kumar, Edward Adelson, and Pulkit Agrawal. Visual dexterity: In-hand dexterous manipulation from depth. *arXiv preprint arXiv:2211.11744*, 2022.

[33] Zhao-Heng Yin, Binghao Huang, Yuzhe Qin, Qifeng Chen, and Xiaolong Wang. Rotating without seeing: Towards in-hand dexterity through touch. *arXiv preprint arXiv:2303.10880*, 2023.

[34] Haozhi Qi, Brent Yi, Sudharshan Suresh, Mike Lambeta, Yi Ma, Roberto Calandra, and Jitendra Malik. General in-hand object rotation with vision and touch. In *Conference on Robot Learning*, pages 2549–2564. PMLR, 2023.

[35] Haozhi Qi, Ashish Kumar, Roberto Calandra, Yi Ma, and Jitendra Malik. In-hand object rotation via rapid motor adaptation. In *Conference on Robot Learning*, pages 1722–1732. PMLR, 2023.

[36] Sudeep Dasari, Abhinav Gupta, and Vikash Kumar. Learning dexterous manipulation from exemplar object trajectories and pre-grasps. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3889–3896. IEEE, 2023.

[37] Gagan Khandate, Siqi Shang, Eric T Chang, Tristan Luca Saidi, Yang Liu, Seth Matthew Dennis, Johnson Adams, and Matei Ciocarlie. Sampling-based exploration for reinforcement learning of dexterous manipulation. *arXiv preprint arXiv:2303.03486*, 2023.

[38] Binghao Huang, Yuanpei Chen, Tianyu Wang, Yuzhe Qin, Yaodong Yang, Nikolay Atanasov, and Xiaolong Wang. Dynamic handover: Throw and catch with bimanual hands. *arXiv preprint arXiv:2309.05655*, 2023.

[39] Toru Lin, Zhao-Heng Yin, Haozhi Qi, Pieter Abbeel, and Jitendra Malik. Twisting lids off with two hands. *arXiv preprint arXiv:2403.02338*, 2024.

[40] Yuanpei Chen, Tianhao Wu, Shengjie Wang, Xidong Feng, Jiechuan Jiang, Zongqing Lu, Stephen McAleer, Hao Dong, Song-Chun Zhu, and Yaodong Yang. Towards human-level bimanual dexterous manipulation with reinforcement learning. *Advances in Neural Information Processing Systems*, 35:5150–5163, 2022.

[41] Kairui Ding, Boyuan Chen, Ruihai Wu, Yuyang Li, Zongzheng Zhang, Huan-ang Gao, Siqi Li, Guyue Zhou, Yixin Zhu, Hao Dong, et al. Preafford: Universal affordance-based pre-grasping for diverse objects and environments. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7278–7285. IEEE, 2024.

[42] Yuanpei Chen, Chen Wang, Li Fei-Fei, and C Karen Liu. Sequential dexterity: Chaining dexterous policies for long-horizon manipulation. *arXiv preprint arXiv:2309.00987*, 2023.

[43] Kevin Zakka, Laura Smith, Nimrod Gileadi, Taylor Howell, Xue Bin Peng, Sumeet Singh, Yuval Tassa, Pete Florence, Andy Zeng, and Pieter Abbeel. Robopianist: A benchmark for high-dimensional robot control. *arXiv preprint arXiv:2304.04150*, 2023.

[44] Nikhil Chavan Dafle, Alberto Rodriguez, Robert Paolini, Bowei Tang, Siddhartha S Srinivasa, Michael Erdmann, Matthew T Mason, Ivan Lundberg, Harald Staab, and Thomas Fuhlbrigge. Extrinsic dexterity: In-hand manipulation with external forces. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1578–1585. IEEE, 2014.

[45] Chengzhong Ma, Houxue Yang, Hanbo Zhang, Zeyang Liu, Chao Zhao, Jian Tang, Xuguang Lan, and Nanning Zheng. Dexdiff: Towards extrinsic dexterity manipulation of ungraspable objects in unrestricted environments. *arXiv preprint arXiv:2409.05493*, 2024.

[46] Simon Stepputtis, Yezhou Yang, and Heni Ben Amor. Extrinsic dexterity through active slip control using deep predictive models. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3180–3185. IEEE, 2018.

[47] Yi Dong, Jinjun Duan, Yangjun Liu, Zhendong Dai, and Poramate Manoonpong. Robotic shoe packaging strategies based on a single soft-gripper system and extrinsic resources. In *2023 International Conference on Advanced Robotics and Mechatronics (ICARM)*, pages 469–475. IEEE, 2023.

[48] Huy Ha and Shuran Song. Flingbot: The unreasonable effectiveness of dynamic manipulation for cloth unfolding. In *Conference on Robot Learning*, pages 24–33. PMLR, 2022.

[49] Tianhao Wu, Yunchong Gan, Mingdong Wu, Jingbo Cheng, Yaodong Yang, Yixin Zhu, and Hao Dong. Unidexfpm: Universal dexterous functional pre-grasp manipulation via diffusion policy. *arXiv preprint arXiv:2403.12421*, 2024.

[50] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[51] George Konidaris and Andrew Barto. Skill chaining: Skill discovery in continuous domains. In *the Multidisciplinary Symposium on Reinforcement Learning, Montreal, Canada*, 2009.

[52] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.

[53] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.

[54] Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield. Foundationpose: Unified 6d pose estimation and tracking of novel objects. *arXiv preprint arXiv:2312.08344*, 2023.

# Appendix

## A. Hyperparameters of the PPO

Table 4 | Hyperparameters of PPO.

| Hyperparameters | Value |
|---|---|
| Num mini-batches | 4 |
| Num opt-epochs | 5 |
| Num episode-length | 8 |
| Hidden size | [1024, 512, 256] |
| Clip range | 0.2 |
| Max grad norm | 1 |
| Learning rate | 5e-4 |
| Discount ($\gamma$) | 0.99 |
| GAE lambda ($\lambda$) | 0.95 |
| Init noise std | 1.0 |
| Desired kl | 0.008 |
| Ent-coef | 0 |

Table 4 shows the hyperparameters of the PPO.

### A.1. Reward Design

Instead of being manually designed for each object, the contact point $c_p \in \mathbb{R}^3$ is defined as a fixed spatial offset from the center point of all objects, maintaining a constant 7 cm distance along the object's width axis as show in Figure 6(a) without manually designed for each object. Given an object with its center point $P_t^{obj} = (x, y, z)$, the contact point is computed as

$$c_p = P_t^{obj} + d \cdot \hat{w} \qquad (3)$$

where $d = 7cm$ represents the fixed offset distance and $\hat{w}$ denotes the unit vector along the object's width dimension.

In Equation 2 , we divide the reward function into three parts: $r_{\text{motion}}$, $r_{\text{pregrasp}}$, and $r_{\text{grasp}}$. $r_{\text{motion}}$ guides the policy toward its ultimate goal and remains active throughout the entire task execution. $r_{\text{pregrasp}}$ encourages the dexterous hand to move towards the object for pre-grasp following the trajectory we expect (first moving towards $c_p$, and then moving towards the object center $P_t^{obj}$). $r_{\text{grasp}}$ is designed to facilitate successful object grasping after pre-grasp.

For $\pi_{\text{wall}}$, we expect that the hand first approaches the object from its side guided by $r_{\text{pregrasp}}$, and then grasp it between the thumb and other fingers guided by $r_{\text{grasp}}$. Thus, we set $P(a) = 1, P(b) = 0$ in the pre-grasp stage and switch to $P(a) = 0, P(b) = 1$ in the grasp stage. The switch occurs when $r_{\text{pregrasp}} < 3cm$ which indicates sufficient proximity between the middle finger $F_t^{f,3}$ and $c_p$. $r_{\text{motion}}$ continuously compute the distance between the object $P_t^{obj}$ and a point $P_t^{target} = P_t^{obj} + [0, 0, 10](cm)$. Notice that this point



(a) $c_p$ position for each task    (c) Importance of $c_p$

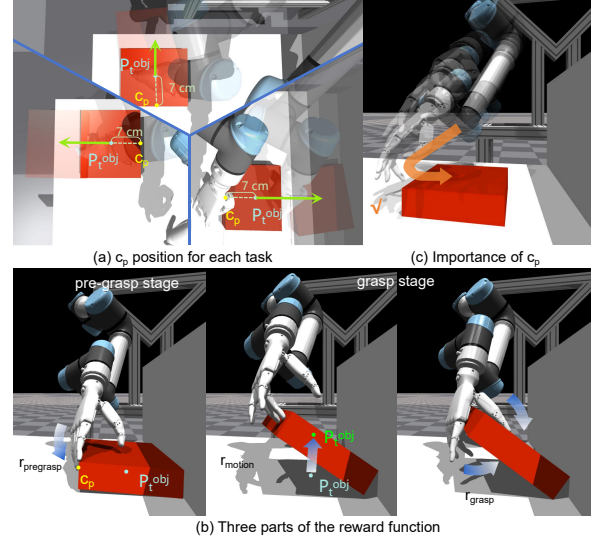pre-grasp stage    grasp stage

(b) Three parts of the reward function

Figure 6 | **Illustration for reward design.** (a) $c_p$ position for each task. (b) Three parts of the reward function. (c) Importance of $c_p$. For $\pi_{\text{wall}}$,

locates above the object which is designed to guide the robot to rotate up the object.

For $\pi_{\text{edge}}$, the object is intentionally positioned to expose a graspable side at the table edge. Therefore, the pre-grasp phase is eliminated ($P(a) \equiv 0$) since moving the finger to $c_p$ becomes unnecessary. Besides, the object translation in the horizontal plane is no longer required, allowing the policy to focus exclusively on vertical finger coordination. The motion reward is formulated as:

$$r_{\text{motion}} = -\sum_{i=1}^{5} |F_t^{f,i} - P_t^{target}| \qquad (4)$$

where target heights are:

$$P_t^{target} = \begin{cases} P_t^{obj} + [0, 0, 0.15]\text{cm}, & \text{for thumb finger} \\ P_t^{obj} - [0, 0, 0.05]\text{cm}, & \text{for other fingers} \end{cases} \qquad (5)$$

We set $P(b) = 1$ if the middle finger moves below the object to activate the grasp stage. Otherwise, we maintain $P(b) = 0$.

For $\pi_{\text{push}}$, stage training and grasping are unnecessary, resulting in $P(a) \equiv 1$ and $P(b) \equiv 0$. $r_{\text{pregrasp}}$ guides the hand to push the object from its side, and $r_{\text{motion}}$ narrows the gap between the object $P_t^{obj}$ and $P_t^{target}$ predicted by our high-level planner.

### A.2. Sim-to-real Details

**Teacher-student distillation.** We collect 1000 demonstration trajectories with the teacher RL policies in simulation for each task. Here we manually design

some rules to remove the unnatural or dangerous behaviors which emerge from exploiting the simulator dynamics but don't transfer well to real-world. We use a transformer-based network to imitate the curated demonstration. The network architecture is as followed:

The network takes as input a sequence of concatenated state observations (dimension: 13) spanning 10 historical frames. An initial feature extraction module processes each frame independently through two linear layers (128 and 512 units respectively), each followed by ReLU activation and layer normalization. The extracted features are augmented with learnable positional encodings to preserve temporal information. The temporal dynamics are modeled through a 3-layer transformer encoder (dmodel=512, nheads=2), where each layer contains: Multi-head self-attention for capturing frame dependencies; Position-wise feed-forward network; Residual connections and layer normalization. Following the transformer encoder, we employ global average pooling across the temporal dimension and process the features through two residual blocks for enhanced representation learning. The final action predictor consists of a carefully designed MLP with progressively decreasing layer widths (256 → 128 units), each followed by ReLU activation, layer normalization, and dropout (p=0.1). The network outputs 8 consecutive action frames (dimension: 12 per frame) through a tanh-activated linear layer, ensuring actions remain within valid bounds.

We supervise the output action $\mathbf{a}_{\text{pred}}$ with negative log product loss with L2 regularization:

$$\mathcal{L}(\mathbf{a}_{\text{pred}}, \mathbf{a}_{\text{gt}}) = -\sum_{i=1}^{N} \log\left(1 - |\mathbf{a}_{\text{pred}}^{(i)} - \mathbf{a}_{\text{gt}}^{(i)}|1\right) + \lambda ||\mathbf{a}_{\text{pred}}||_2^2 \tag{6}$$

**Digital twin.** Before leveraging the teacher-student

Table 5 | Results for the real-world experiments

| | box-w1 | box-w2 | box-w3 | bag-w1 | container | handbag |
|---|---|---|---|---|---|---|
| Wall | 6/10 | 8/10 | 7/10 | 9/10 | 9/10 | 9/10 |
| | box-e1 | box-e2 | bag-e1 | bag-e2 | plate | handbag |
| Edge | 10/10 | 7/10 | 10/10 | 8/10 | 5/10 | 9/10 |

distillation, we achieve zero-shot sim-to-real transfer by implementing a digital twin framework that bridges our simulation policy with the real dexterous arm-hand system. The framework operates through two parallel threads that enable real-time asynchronous communication between simulation and real-world environments. The real-world thread continuously collects observations, including arm-hand proprioception and object pose information. Meanwhile, the

simulation thread processes these observations to generate control actions, executes them in simulation and uses the resulting joint angles as target joint angles for PD control in the real robot system. The simulation environment is continuously synchronized with real-world by updating robot joint angles and object poses from real-world observations. Through the constant synchronization of simulation and real-world, we evaluate our RL polices following the similar setup as mentioned in Subsection 5.4. The results in Table 5 shows that our digital twin framework achieves robust performance on real-world objects, with success rates exceeding 80% for most objects.