

PP-Tac: Paper Picking Using Omnidirectional Tactile Feedback in Dexterous Robotic Hands

Author Names Omitted for Anonymous Review. Paper-ID 143

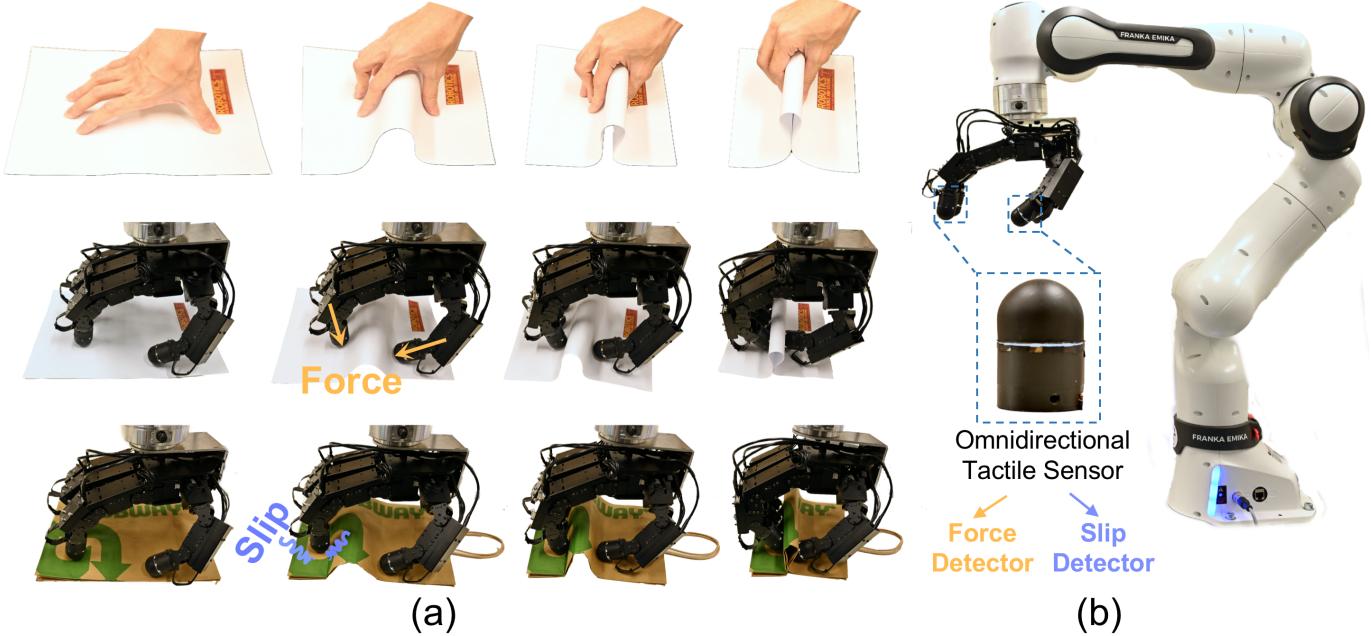


Fig. 1: **Overview of the PP-Tac Robotic System.** The system leverages omnidirectional tactile feedback in a dexterous robotic hand to pick up thin, deformable paper-like objects. (a) Hand motion generated by force and slip feedback, inspired by human paper-picking motions involving sliding and pinching. (b) Hardware setup includes a robotic arm, a dexterous hand, and four fingertip-mounted tactile sensors capable of simultaneously detecting force and slip events.

Abstract—Robots are increasingly envisioned as human companions, assisting with everyday tasks that often involve manipulating deformable objects. While recent advancements in robotic hardware and embodied AI have expanded their capabilities, current systems still struggle with handling thin, flat, deformable objects like paper and fabric due to limitations in motion planning and perception. This paper introduces PP-Tac, a robotic system for picking up paper-like objects. PP-Tac features a multi-fingered robotic hand equipped with high-resolution tactile sensors, providing omnidirectional feedback for slip detection and precise friction control. Additionally, we propose a grasp trajectory synthesis pipeline that generates a dataset of paper-like object grasping motions and trains a diffusion-based motion generator, which is then implemented on a physical hand-arm platform for evaluation. Experiments demonstrate PP-Tac’s effectiveness in grasping paper-like objects of varying stiffness (e.g., cloth and paper), achieving a success rate of 87.5%. By leveraging tactile feedback, PP-Tac adapts to varying surfaces beneath the objects with robustness. This study is the first to explore grasping thin, deformable objects using a dexterous robotic hand with tactile feedback. These advancements pave the way for broader applications in domestic, industrial, and logistical settings, where precise handling of paper-like objects is essential.

I. INTRODUCTION

Robots are increasingly popular as assistive agents in everyday life, particularly within household environments [34].

These robots are designed to perform various domestic tasks, often involving the grasp of thin, deformable objects such as paper and fabric [47]. For instance, clothes-folding tasks [25] require high dexterity and adaptability to accommodate variations in fabric size, texture, and stiffness, while document organization tasks [2] demand precise picking capabilities for diverse paper types and form factors. Beyond domestic settings, the ability to handle deformable objects is essential in industrial and logistical applications, such as fabricating fabrics [4] and packing objects using plastic bags and cardboard [11].

Despite their significance, picking up paper-like objects remains challenging in robotics [47]. In particular, the main challenges are three-fold: 1) Vision systems, commonly used for manipulation, struggle to perceive contact information during interactions with deformable objects due to limited sensing modalities and occlusion, resulting in a lack of necessary feedback for motion planning [26]; 2) Their thin, stiff characteristics often result in flat shapes, hindering the synthesis of stable grasps using conventional methods due to insufficient contact points [8]. 3) The appearance of such objects exhibits high variability, as their shape undergoes continuous and unpredictable deformation during manipulation. These dynamic shape variations significantly impair the generalizability of

vision-based methods.

In contrast, humans excel at picking up paper-like objects by leveraging coordinated multi-fingered motion and tactile sensing. As shown in Fig. 1(a), the process typically begins with establishing contact using fingers, followed by sliding motions to deform the material and enable a stable pinched grasp. Such success stems from the coordination of multiple fingers, which generates the necessary friction to deform the object and utilizes sufficient finger Degree of Freedoms (DoFs) to adaptively establish stable contact points for a stable grasp. Additionally, tactile sensing complements visual feedback, allowing humans to perceive the object’s deformation and apply appropriate forces by detecting friction and slip. These tactile cues facilitate real-time adjustments, ensuring the successful execution of the picking-up action.

Inspired by human strategies, this paper introduces a robotic system, *PP-Tac* (Paper-like object Picking using omnidirectional Tactile feedback), designed for dexterous robotic hands. The system comprises two key components: **A dexterous robotic hand with omnidirectional and high-resolution Vision-Based Tactile Sensors (VBTS)**. These fingertip-mounted tactile sensors provide real-time feedback on contact status during grasping and feature an omnidirectional sensing area with a high-framerate monochrome camera, enabling faster response times and simpler calibration compared to RGB-based systems. An illustration of the system is shown in Fig. 1(b). In addition to the tactile sensor, this paper also presents **A diffusion-based motion generation policy (PP-Tac policy)** that imitates human picking-up skills. The proposed method first employs efficient trajectory optimization to generate expert data replicating human sliding and pinching motions. To generalize this approach to diverse deformable objects and uneven surfaces, a diffusion policy is subsequently trained using these trajectories, leveraging proprioceptive data and tactile feedback for adaptive control of the dexterous robotic hand.

In comprehensive real-world experiments, the proposed PP-Tac achieved an overall success rate of 87.5% in grasping everyday thin and deformable paper-like objects, such as plastic bags, paper bags, and silk towels on flat surfaces. Fig. 1(a) illustrates examples of our arm-hand system successfully picking up paper-like objects. The PP-Tac also demonstrates significant adaptability in picking up paper-like objects on various uneven surfaces. Additionally, an ablation study further validates the contributions of each system component, highlighting the critical role of VBTS feedback and motion generation policies in achieving robust performance.

To the best of our knowledge, this work represents the first demonstration of deformable object picking using a dexterous hand equipped with VBTS. Overall, our contributions include:

- 1) We propose a new omnidirectional tactile sensor that is easy to fabricate, calibrate, and deploy at scale.
- 2) We assemble a fully actuated dexterous robotic hand integrated with VBTS into each fingertip to enable real-time contact feedback.
- 3) We introduce *PP-Tac policy*, a diffusion policy for picking

up paper-like objects that demonstrate robust generalization across diverse materials and surfaces.

- 4) We provide the implementation and systematic experiments of the proposed algorithms on the physical device. Upon the paper’s acceptance, both hardware and code for *PP-Tac system* will be open-sourced to support further research and community development.

II. RELATED WORK

A. Deformable Objects Manipulation

Deformable Object Manipulation (DOM) involves handling soft objects that alter shape during interaction—a ubiquitous yet challenging task in robotics. Challenges arise from uncertainties in perception and complex soft-body dynamics [14, 3]. Early approaches relied on visual perception for state estimation [47], enabling tasks like rope-handling [30, 44] and cloth-folding [37, 25]. However, vision-based methods often underperform in real-world DOM due to varying object appearance, limited physical property perception, occlusions [23, 5], and variable lighting conditions [43, 20]. These challenges hinder the development of scalable vision-based DOM solutions for diverse environments.

Tactile sensing, particularly Vision-Based Tactile Sensors (VBTS), has demonstrated significant potential for DOM tasks compared to visual perception [47]. VBTS excel in tasks such as object shape reconstruction [31, 9, 28], localization [18, 24, 6], and slip detection [39, 12], leveraging their high-resolution tactile feedback. Prior work has explored VBTS for deformable object manipulation [35], but existing implementations rely on gripper-mounted sensors, which lack the dexterity of multi-fingered hands due to limited DoF. Our experiments reveal that gripper-based approaches struggle with thin, deformable objects or those on non-flat surfaces, highlighting the need for integrating dexterous robotic hands with VBTS for robust manipulation [19].

B. Dexterous Robotic Hand with Tactile Sensing

Current dexterous hands are often equipped with tactile sensing capabilities to enhance dexterity. Commonly used tactile sensors typically incorporate mechanisms such as capacitive [29], piezoresistive [17], or magnetic-based [13] technologies. These designs can be fabricated in various shapes and sizes, allowing them to conform to the form factor of different robotic fingers. However, the sensing principles behind these technologies limit their spatial resolution and robustness under varying environmental conditions. Although efforts have been made to develop VBTS for curved shapes [9, 41, 42, 21], these designs are not yet commercially available and remain challenging to deploy at scale in hand configurations. To address this, we propose an omnidirectional VBTS that is structurally simple, compact, easy to fabricate, and straightforward to calibrate. Its monochrome sensing principle inherently minimizes required data bandwidth while enabling precise curved surface reconstruction, making it well-suited for large-scale deployment on dexterous robotic hand systems. Additionally, current robotic hands equipped with VBTS have been used

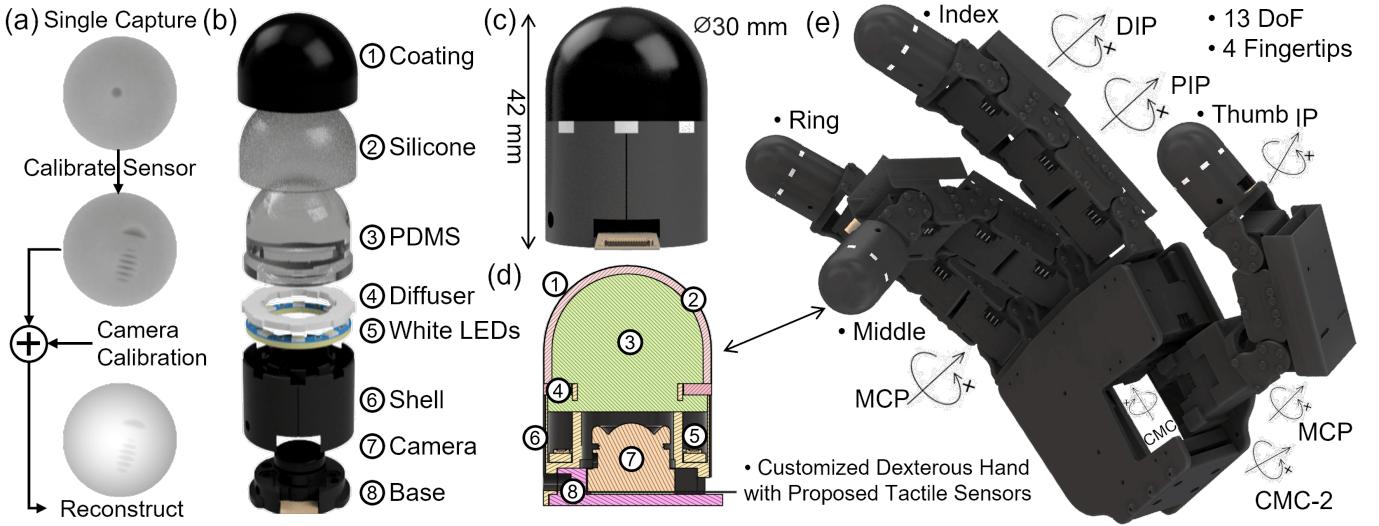


Fig. 2: The hardware design of the omnidirectional VBTS and its integration into the four-fingered dexterous robotic hand system. (a) illustrates the pipeline of depth reconstruction. (b) illustrates the exploded view of the sensor, detailing each component. (c) shows the dimensions of the sensor. (d) shows the schematic design. (e) illustrates the robotic hand equipped with four sensors on its distal joint.

in grasping and in-hand orientation tasks. For instance, Do *et al.* uses DenseTact, attached to an Allegro Hand, to grasp and manipulate small screws [10]. Qi *et al.* integrates fingertip VBTS [32] and DIGIT [38] on an Allegro Hand to enable in-hand object rotations. To the best of our knowledge, no previous research has been conducted on VBTS sensorized dexterous hands for picking up thin and deformable objects, such as paper-like materials.

III. HARDWARE DESIGN

To provide sufficient dexterity to address the challenges of paper-picking tasks, we designed and fabricated a set of finger-shaped VBTS, which are then integrated into Allegro Hand [33] through customization.

A. Fingertip-shaped Tactile Sensing

The design of the fingertip-shaped tactile sensor is guided by five key principles to ensure effective manipulation:

- **Round shape:** The hemispherical design enables omnidirectional tactile perception.
- **High resolution:** High spatial resolution enables accurate force and slip detection during the picking-up process.
- **Ease of fabrication & low-cost:** The components of the tactile sensor are either off-the-shelf or easy to fabricate, with a cost of around \$60.
- **Efficient calibration:** The monochrome sensing principle simplifies lighting control and reduces manual effort in image capture for calibration, making it particularly suitable for large-scale deployment on multi-fingered robotic hands.
- **Efficient data transmission:** The monochrome camera produces lightweight data per frame, facilitating high-speed data transmission between systems.

Based on these principles, the sensor design and its integration into the dexterous robotic hand is illustrated in Fig. 2. Next, we detail each component and the calibration process.

1) Contact and Illumination Module: The core of the sensor is a contact module (elastomer) with a uniformly illuminated, deformable sensitive surface that maintains structural rigidity during contact. Inspired by the monochrome sensing principle [27], where brightness changes indicate deformation, we developed a hemispherical structure comprising a white LED ring, a stiff transparent internal skeleton, a soft semi-transparent perception layer, and a thin opaque protective layer that achieves the desired optical characteristics.

The LED ring (LUXEON 2835 4000K SMD LED) and a diffuser (double-sided frosted diffuser sheet) are first installed within the sensor shell. The skeleton is then manufactured from PDMS (Dow Corning Sylgard 184 with Shore hardness 50 A) using a two-piece molding technique. The mixture (base: catalyst = 10: 1) is degassed and poured into the mold, and cured for 24 hours at room temperature. The perception layer is then manufactured similarly, using semitransparent silicone (Smooth-On Ecoflex with Shore hardness 00-10), and the layer is peeled off after 4 hours. Finally, a silicone coating (Smooth-On Psycho Paint) is airbrushed onto the perception layer to form the opaque protective layer. The entire manufacturing process takes within 3 days, facilitating large-scale deployment.

2) Camera Module: A micro black-and-white CMOS camera with a wide 160° lens is used to capture single-channel brightness data. The camera operates up to 120Hz with a resolution of 640 × 480 and a latency of approximately 100ms.

3) Calibration: The uniform optical properties of the elastomer and illumination module (with a capture standard deviation as low as 6) enable the 3D geometry of the rounded sensor to be computed from single-channel pixel intensity in simply two steps using only 30 captures. First, given the known intrinsic parameters K , camera calibration is performed using 29 captures in an indentation-based setup to estimate the extrinsic parameters of rotation matrix A and translation

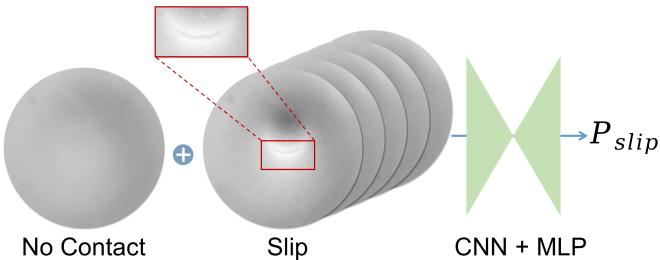


Fig. 3: **Slip Detection.** The left tactile image shows no contact, while the middle tactile image highlights wrinkle features during slip. The network computes the probability of slip P_{slip} using the no-contact tactile image and five sequential tactile images.

vector b , as well as the sensor surface reference projection D . Next, the depth mapping function M is calibrated by capturing a single image of a ball of known size pressed onto the sensor [27]. The complete mapping function from the pixel coordinates (u, v) to the sensor coordinates (x, y, z) can be expressed as:

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = A^{-1} \left((D(u, v) - M(I_\Delta(u, v))) K^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} - b \right), \quad (1)$$

which allows grayscale intensity images to be transformed to a depth map expressed in the sensor coordinates. A detailed explanation of camera calibration is provided in [Appendix A](#). Reconstruction results and qualitative analysis are presented in [Section VI-B](#).

4) *Contact Force Estimation & Slip Detection:* Our sensors are capable of detecting both contact forces and slip events. The contact force, modeled by elasticity theory, is proportional to the deformation depth and can be expressed as a function of deformation depth. Furthermore, the slip between the sensor and the object surface is detected using a lightweight neural network, as described in [Fig. 3](#). The network takes the previous five frames as inputs, extracts the features via a convolutional neural network (CNN), and outputs the slip probability P_{slip} through a multilayer perception network (MLP). To train this network, we collected approximately 20 minutes of data from the four tactile sensors. When the threshold of P_{slip} is set to 0.75, our evaluation shows that the system achieves a detection accuracy of 86%.

B. Robotic Hand System

We integrated the proposed omnidirectional tactile sensors into a fully actuated dexterous robotic hand. These tactile sensors are mounted at the distal end of each fingertip, facilitating contact characterization in the following paper-picking tasks. We designed and fabricated the robotic hand featuring 13 controllable DoFs, including the DIP, PIP, and MCP joints for the index, middle, and ring fingers, as well as the CMC, CMC-2, MCP, and IP joints for the thumb. The robotic hand is driven by Dynamixel XC330-M288-T motors, which are all multiplexed through a U2D2 Hub. For each tactile sensor, it communicates with the PC via a USB interface. The entire assembly is mounted on a Franka

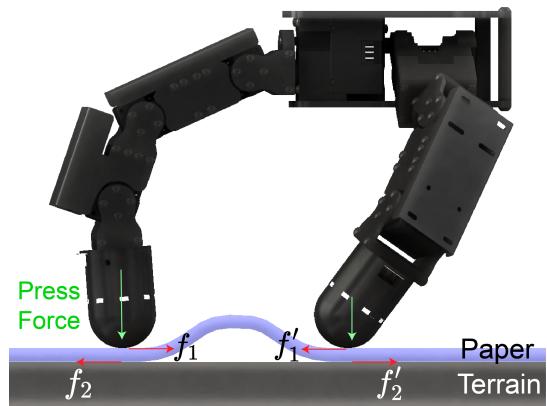


Fig. 4: **Force analysis during grasping flat objects.** The grasping process is made possible by the following forces: (1) the contact normal force exerted by the sensor on the object. (2) the static friction force (f_1, f'_1) between fingers and the object, (3) a dynamic friction force (f_2, f'_2) between the object and the terrain. When the static friction (f_1, f'_1) exceeds the critical buckling resistance of the paper, the sheet deforms, creating a stable pinch region that facilitates successful grasping.

Research 3, a 7-DoF robotic arm, which communicates with the PC via a high-speed Ethernet connection.

IV. PAPER-LIKE OBJECT PICKING PROBLEM

Next, we aim to address the challenge of grasping thin, deformable paper-like objects from flat surfaces. This appears as a commonly seen scenario in everyday tasks, such as organizing scattered document pages or retrieving napkins from dining plates. Although creases or irregularities in the material can sometimes provide grasping points, a particularly challenging scenario arises when the object is extremely flat and lacks discernible edges or salient grasping features. This research introduces a novel approach to tackle this paper-picking problem that was previously unexplored.

Motivated by the human strategy for grasping flat objects, our work is based on a biomimetic grasping pose optimized for paper picking, as illustrated in [Fig. 4](#). By applying sufficient inward force, the robotic fingers can induce buckling of the material against the supporting surface. This buckling effect dynamically generates a pinchable region, enabling subsequent grasp execution.

During buckling, the distance between contact points beneath the fingers decreases. When this reduction rate matches the fingertips' closure speed (*i.e.*, no relative motion between fingertips and material), two frictional forces govern the system: static friction (f_1, f'_1) between the fingers and material, and dynamic friction (f_2, f'_2) between the material and the supporting surface. Their magnitudes depend on the applied normal force and the respective coefficients of friction.

In particular, the above analysis assumes that the static friction between robotic fingers and the material exceeds both the maximum static friction at the material-terrain interface and the critical buckling resistance of the material. This framework can also be extended to scenarios with uneven supporting surfaces. Without loss of generality, we assume that height variations in the terrain are less than 3 cm.

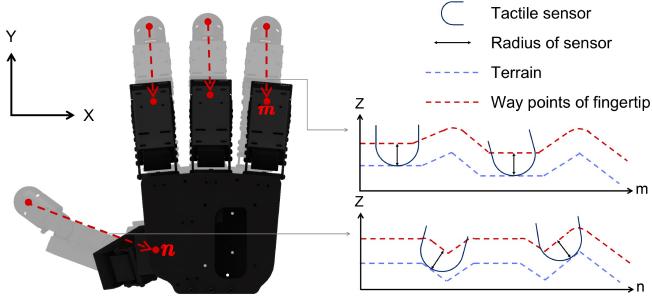


Fig. 5: **Fingertip trajectories from data synthesis.** Trajectories ensure fingertip sliding along the terrain surface. Adjusting the distance between waypoints and terrain affects sensor deformation. The right figure projects trajectories of two fingers onto the m - z and n - z planes, where m and n are straight-line projections of fingertip trajectories on the palm-aligned x - y plane, and the z -axis extends outward from the hand.

V. POLICY LEARNING FOR PAPER-PICKING

Manipulating paper-like objects with visual perception remains challenging due to difficulties in detecting thickness and textural variability. To address this, we propose a vision-independent tactile-based approach. The core idea leverages tactile feedback to maintain contact conditions (as defined in [Section IV](#)), facilitating the creation of a buckling region for successful grasping. We implement this through the *PP-Tac policy*, developed in two stages: 1) Trajectory Optimization: Generate a dataset of grasping motions using trajectory optimization. 2) Diffusion Policy Training: Train a policy on this dataset to infer motions from tactile feedback and proprioceptive states, ensuring generalization to real-world robotic systems.

A. Grasp Motion Dataset Synthesis

We synthesize grasping motions through trajectory optimization in simulation, avoiding the need for complex teleoperation devices. While reinforcement learning (RL) offers an alternative, it requires soft-body simulation to model deformable object dynamics and VBTS elastomer behavior, often necessitating additional real-to-sim procedures for fidelity. In contrast, our approach uses rigid-body dynamics and transfers directly to real robots, as validated experimentally. The grasping process begins by establishing fingertip contact with the object's surface (see [Appendix B](#) for implementation details). Once contact is achieved, the fingers gradually close to complete the grasp. Each finger follows an independent trajectory on the object's surface, with normal forces adjusted to maintain contact ([Figs. 4 and 5](#)).

In simulation, the ground-truth shape of the terrain is known, enabling the determination of all finger joint values

and arm poses through the following optimization problem:

$$\hat{\gamma} = \arg \min_{\gamma} (L_{ee} + L_{\Delta} + L_{RT}), \quad (2)$$

$$L_{ee} = w_{ee} \text{MSE}(\mathbf{fk}(\gamma), ee_{target}), \quad (3)$$

$$L_{\Delta} = w_{\Delta} \text{MSE}(\hat{\gamma}, \gamma), \quad (4)$$

$$L_{R,p_{wrist}} = w_{R,p_{wrist}} \text{MSE} \left((\hat{R}^{1:N_{data}}, \hat{p}_{wrist}^{1:N_{data}}), (R^{1:N_{data}}, p_{wrist}^{1:N_{data}}) \right), \quad (5)$$

where γ is the optimization variables consisting of hand joint angles $q^{1:N_{data}}$; $R^{1:N_{data}}$ is the rotation matrix of wrist(end effector of arm) rotation, and $p_{wrist}^{1:N_{data}}$ is the wrist translation along the z -axis in world coordination; N_{data} is the sequence length. The forward kinematics \mathbf{fk} computes the four fingertips' trajectories, and ee_{target} represents the target fingertips' trajectories. The objective function minimizes the mean squared error (MSE) between the fingertip positions and their targets, while L_{Δ} regularizes the motion to remain close to the initial pose. Additionally, $L_{R,p_{wrist}}$ minimizes wrist movement, ensuring the arm stays within its workspace.

To account for varying material properties, sliding trajectories are adjusted based on fingertip contact forces. This is achieved by synthesizing motions that deform the tactile sensor's elastomer layer to different extents. The maximum deformation reading d_{tac} is proportional to the applied pressure, governed by the elastomer's Young's modulus. Thus, contact force F is modulated by controlling d_{tac} via position control. Notably, the exact relationship between d_{tac} and F is not explicitly modeled, as precise force values are unnecessary for the algorithm. This approach leverages rigid-body dynamics to control contact forces efficiently, avoiding complex deformable dynamics calculations. By adjusting the distance between the finger joint and the terrain, we can obtain trajectories with varying degrees of deformation. When the distance between the finger joint and the terrain is equal to the sensor's radius as shown [Fig. 5](#), the finger just makes contact with the terrain and d_{tac} just equals to 0. Using this method, we generated a dataset of 500,000 grasp samples, each comprising a sequence of $N_{data} = 100$ frames.

B. PP-Tac Policy

Once the dataset is prepared, we employ a diffusion policy to jointly control the hand and arm, enabling adaptation to varying terrain shapes and contact force conditions. We adopt a Denoising Diffusion Probabilistic Model (DDPM) framework [15, 16, 7, 36], which predicts future actions (N_{pred} steps of x^{pred}) conditioned on historical states (N_{prefix} steps of x^{prefix}). The state variables include:

$$(\mathbf{p}_j, \dot{\mathbf{p}}_j, \mathbf{q}, \dot{\mathbf{q}}, R, \Omega, p_{wrist}, \dot{p}_{wrist}, \mathbf{d}_{tac})$$

where $\mathbf{p}_j \in \mathbb{R}^{17 \times 3}$ is hand joints' position in world coordinate, $\dot{\mathbf{p}}_j \in \mathbb{R}^{17 \times 3}$ is the linear velocity of the hand joints relative to each parent frame, $\mathbf{q} \in \mathbb{R}^{13}$ is the rotation angle of controllable hand joints, $\dot{\mathbf{q}} \in \mathbb{R}^{13}$ is the angular velocity of controllable hand joints, $R \in \mathbb{R}^6$ is 6D rotation (represented as two-row

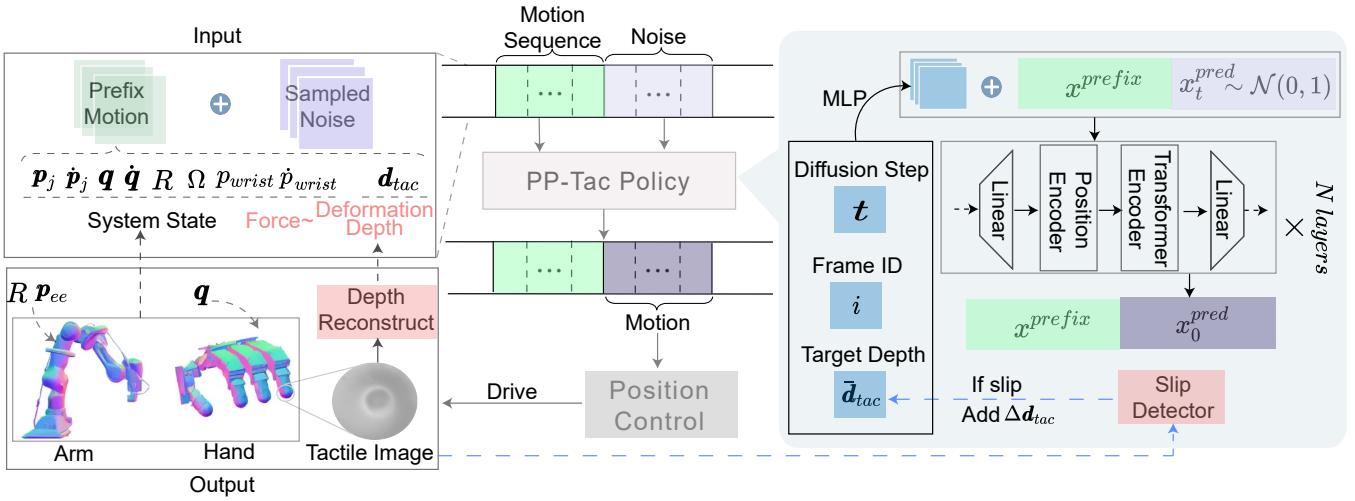


Fig. 6: **Inference pipeline of the proposed PP-Tac policy.** Conditioned on robot proprioception and the target force that needs to be exerted, PP-Tac can infer the action of the next steps. If slip is detected between the finger and the flat object underneath, an incremental amount of force will be exerted by the finger.

vectors of rotational matrix, which is from [46]) of wrist(end effector of arm), $\Omega \in \mathbb{R}^6$ represents the angular velocity of wrist rotation, $p_{wrist} \in \mathbb{R}$ is the wrist's height along arm's z -axis, $\dot{p}_{wrist} \in \mathbb{R}$ is the linear velocity of p_{wrist} , $d_{tac} \in \mathbb{R}^4$ represents the deformation depth readings from four fingertip tactile sensors. Table II summarizes the notations used in this paper. The total state dimension is $D = 142$. Such a high-dimensional and over-parameterized input allows the network to extract more robust and expressive latent features for the diffusion policy.

The pipeline is illustrated in Fig. 6 and Fig. 6 (right) illustrates a single denoising diffusion step. We apply an encoder-only transformer to predict future robot motion x_0^{pred} given prefix motion x^{prefix} , diffused future motion x_t^{pred} , diffusion step t , current frame index i , and target deformation depth \bar{d}_{tac} . The input sequence is encoded into a latent vector of dimension $\mathbb{R}^{(1+N_{\text{prefix}}+N_{\text{pred}}) \times D}$, comprising: 1) A latent vector of D -dimensional features representing t , i , and \bar{d}_{tac} , extracted using three 3-layer MLP networks respectively. 2) $N_{\text{prefix}} \times D$ dimensions corresponding to the prefix states of N_{prefix} time steps. 3) $N_{\text{pred}} \times D$ dimensions for the predicted states of N_{pred} time steps. Instead of predicting ϵ_t as formulated by [16], we follow [40] to predict the state sequence itself \hat{x}_0^{pred} . Predicting \hat{x}_0^{pred} is found to produce better results for the state sequence which contains motion data, and enables us to apply a target loss as geometric loss explicitly as each denoising step as following:

$$L = \|\hat{x}_0^{\text{pred}} - x_0^{\text{pred}}\|_2^2 + \lambda_{\text{consist}} L_{\text{consist}}, \quad (6)$$

$$L_{\text{consist}} = \|\mathbf{fk}(q_0^{\text{pred}}) - J_0^{\text{pred}}\|_2^2 \quad (7)$$

where L_{consist} enforces consistency between joint angles and positions, and λ_{consist} is a weight hyper-parameter.

During inference, we set $t = 1000$ and the diffused $x_{1000}^{\text{pred}} \sim \mathcal{N}(0, I)$ and iteratively denoise it to produce x_0^{pred} . To ensure

real-time performance, we reduce denoising steps to 10 and set $N_{\text{pred}} = N_{\text{prefix}} = 5$, achieving motion generation in 11 ms on an RTX4090 GPU. The predicted q controls the hand, while R and p_{wrist} control the arm.

During grasping, preventing slip between the object and the fingertips is essential to maximize material deformation. To achieve this, a fingertip contact force controller is introduced, which adjusts the fingertip's deformation depth d_{tac} . If slip is detected by the tactile sensors, we increase the desired deformation depth by a small increment Δd_{tac} .

To deploy diffusion policy to real robots, we also need to tackle the domain gap between the real world and simulation. This is achieved by introducing four distinct ways to incorporate disturbances into x^{prefix} during training.

- Add random Gaussian noise to γ to simulate various control errors that may occur in real-world situations.
- Add Gaussian noise to the first frame and gradually amplify it in subsequent frames, simulating the fingers moving across a rising or descending terrain.
- Randomly choose from 2 to N_{prefix} temporal consistent frames to be static, simulating fingers getting stuck due to excessive pressure on complex terrain. And d_{tac} is set to its maximum threshold. The reason for adding the index of the frame into the input is also to avoid issues caused by the fingers getting stuck.

VI. EXPERIMENTS

In this section, we present comprehensive experiments to evaluate our proposed PP-Tac pipeline. First, we detail the implementation of our algorithm (Section VI-A). Next, we show the quantitative and qualitative results of the depth reconstruction of our VBTS (Section VI-B). Then, we perform systematic comparisons of our system on different flat materials and supporting terrains (Section VI-C). We also compare our system with various manipulators to highlight its advantages and limitations (Section VI-D). Last, ablation

studies are conducted to examine the influence of parameters and the necessary training steps (Section VI-E).

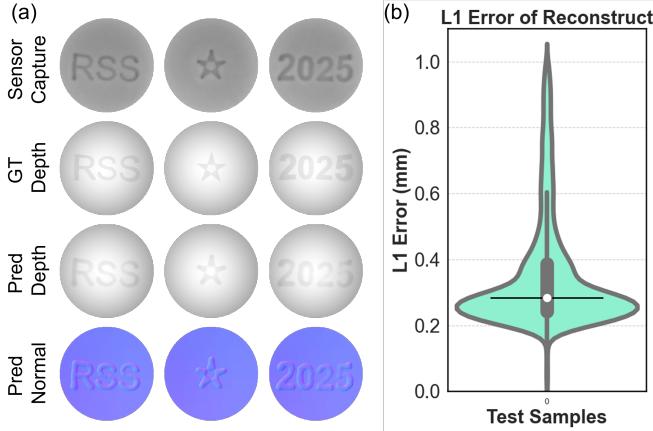


Fig. 7: Reconstruction results. (a) Gallery of reconstructed depth and normal maps from tactile images. (b) Depth reconstruction error of indentation test.

A. Implementation Details

For reproducibility, we provide the implementation details of the PP-Tac algorithm. Our diffusion policy is implemented as a four-layer Transformer encoder with a latent dimension of 512 and four attention heads. We split the each synthesized data sequence into subsequences of length 10 for diffusion process, and train the model for approximately 600,000 iterations on a single RTX 4090. During training, the diffusion step t is uniformly sampled from 0 to 1000. During inference, an acceleration technique is applied as follows. First, t is initialized to 1000 and directly denoised to x_0^{pred} . Subsequently, noise is added to the $t = 1000 - 100N_i$ level and denoised again to x_0^{pred} , where N_i is the inference step number. Thus, the entire inference process consists of 10 steps.

For terrain generation, we model the terrain beneath each finger as a cubic spline with a trajectory length of 100. Control points are placed at intervals of 25 along the trajectory, resulting in a total of 5 control points. To simulate ramps, the height of each control point is randomized by sampling uniformly within the range of $[0, 3]$ cm.

B. Depth Reconstruction of VBTS

To evaluate the performance of the tactile sensor in depth reconstruction, the sensor surface is pressed with three indenters, each with the text content “RSS”, “★” and “2025”. The qualitative results of the sensor output are shown in Fig. 7, which demonstrates the raw captured image from the sensor, the ground truth depth maps, predicted depth maps, and the corresponding calculated normal maps, respectively. These results demonstrate that the sensor can fully reconstruct fine surface details.

We quantify the reconstruction error using a violin plot, leveraging ground truth indentation information obtained from 3D-printed hemispherical shape indicators containing various testing indenters. We collected 215 testing configurations,

each with paired sensor outputs and ground truth reprojection images. The sensor achieves a mean absolute error (L1 error) reconstruction loss of 0.35 mm, and a median loss of 0.28 mm, with 60% of reconstruction losses below 0.3 mm. In terms of computational speed, the depth mapping process takes less than 10 ms, ensuring real-time performance for robotic applications.

C. Evaluation of PP-Tac Policy on Materials and Terrains

We conducted experiments to evaluate the system’s ability to handle flat objects under varying conditions. The qualitative and quantitative results are shown in Fig. 8 and Fig. 9 respectively. Fig. 8 shows the typical successful grasp cases, highlighting that our hardware and PP-Tac algorithm can successfully handle flat objects placed above both the flat and uneven object surface. During the grasping process, the fingertip first contacts the material, followed by a gradual finger closure that buckles the material and creates pinchable regions. Finally, the object is pinched and lifted.

Fig. 9 provides quantitative analysis of the success rate with respect to the object material and the complexity of the terrain beneath. To facilitate this analysis, we conducted experiments using four flat objects in daily life: paper, plastic bag, cloth, and kraft paper bag, each of which presents unique challenges. The paper is extremely flat with no detectable hold points. Plastic bags, commonly encountered in daily life, are difficult to locate using conventional visual pipelines because of their transparency. The cloth is thick and highly deformable, while the kraft paper bags are stiff and have a multilayered structure. To assess the system’s robustness, we also varied the terrain beneath the objects. The four types of terrain used include: a flat plane, a slope (10 degrees), a plane with a 2 cm thick book randomly placed on it, and an uneven terrain with random curvatures. The terrain shapes are shown in Fig. 9.

For statistical significance, we performed 20 grasping attempts for each combination of terrain and object. From results in Fig. 9, cloth and plastic bags are relatively easy to grasp due to their low stiffness, which allows them to buckle more easily under force. In contrast, paper and kraft paper bags are being stiffer and resist buckling, leading to lower success rates.

The terrain beneath the object also significantly impacts grasp success. On flat terrains, such as a plane or a tilted slope, success rates for paper, plastic bags, and cloth were relatively high. This suggests that flat surfaces usually generate consistent frictional forces essential for a successful grasp. However, this advantage diminishes for stiffer flat objects, such as kraft paper bags. These stiff flat objects usually lack of initial buckling when placed on a flat surface, making it more challenging to form reliable grasp points afterward.

For uneven surfaces, the success rates varied according to the shape of the terrain. When a book was placed underneath the flat object, all objects maintained high success rates. These results can be attributed to the edge of the book and the partial void space created beneath the material, which made it easier for the materials to buckle and separate with the terrain. In contrast, when the terrain was highly irregular, the success

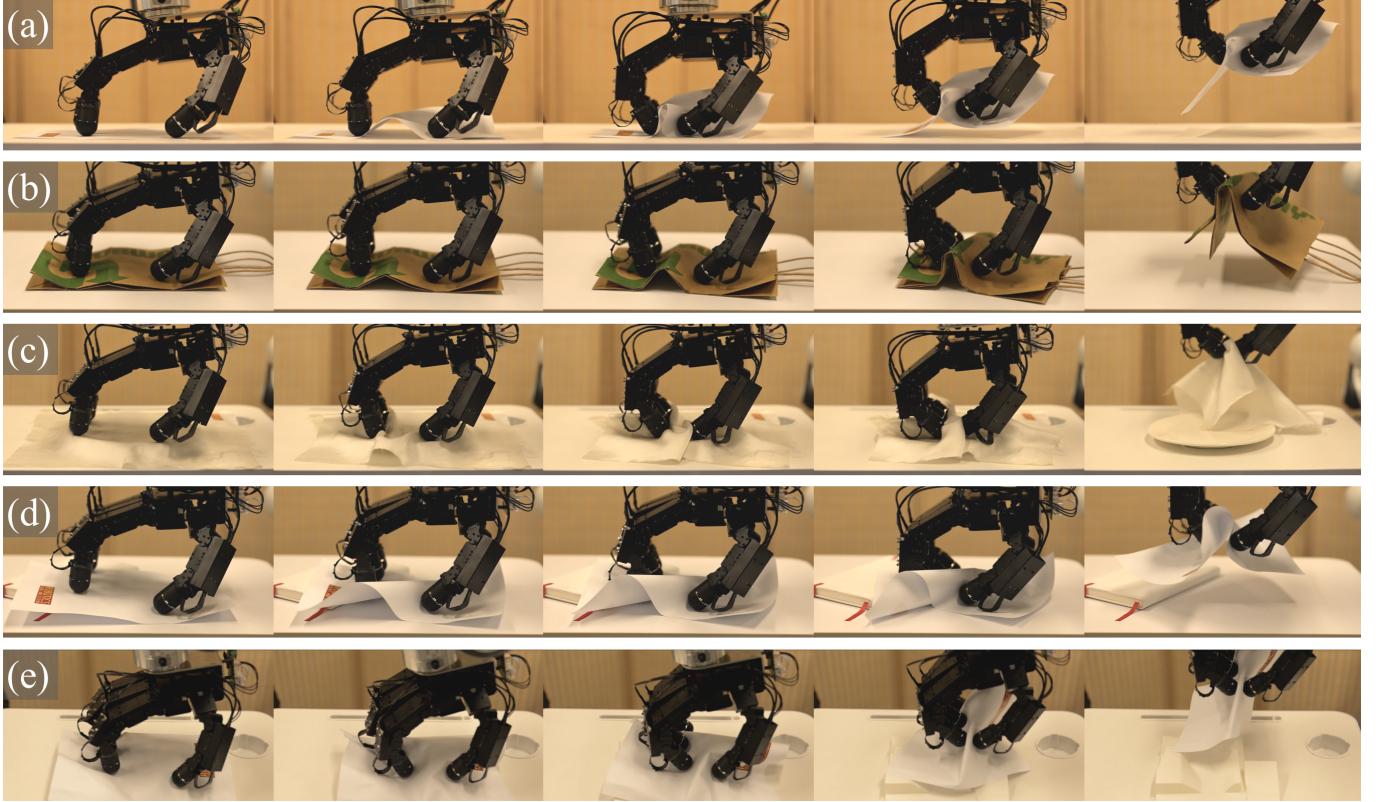


Fig. 8: Gallery of Grasping Different Objects in Real-World Evaluations. This figure demonstrates successful grasps of five flat objects on four different types of terrains, highlighting the effectiveness of our hardware and the PP-Tac algorithm. (a) A paper sheet on a flat desktop. (b) A stiff kraft paper bag on a flat desktop. (c) A soft napkin on a plate. (d) A paper sheet on a randomly arranged book. (e) Paper sheet on a random terrain. These evaluations showcase the robustness and adaptability of our approach.

rate dropped for all objects. This is likely due to the challenges added to our force controllers, which increased the likelihood of the fingers slipping away from the material.

TABLE I: Experimental Results for Varying Paper Quantities: The system’s performance was evaluated on paper materials with different buckling strengths, achieved by bonding 1, 3, 5, and 7 layers of paper with adhesive. For each configuration, 20 trials of grasps were conducted. The average number of slip events detected (No. Slip) and the final success rate (Succ. Rate) were recorded.

Paper Layers	No. Slip	Succ. Rate (%)
1	0.2	90
3	2.9	75
5	13.3	30
7	18.2	5

D. Comparison with Other System Configurations

To assess whether PP-Tac’s system setup leveraging dexterous hand and tactile sensors can offer advantages, systematic comparisons with other robot configurations were conducted. Here, we constructed baselines that include the following two robot configurations. First, robots conventionally use bi-finger grippers rather than dexterous hands. Thus, one strong baseline involved is a bi-finger gripper controlled via human

teleoperation, with a camera mounted on the wrist to provide an egocentric view. To ensure fairness, each trial allowed only one grasp attempt.

In addition to the bi-finger gripper baseline, we compared the dexterous hand’s grasping performance with and without tactile feedback (*i.e.*, open-loop control). For baseline setup, we pre-generated trajectories using the ground truth shape of the terrain and then replayed these trajectories rather than using the PP-Tac policy. Note that this trajectory-replay setting is unattainable in scenarios with high variations, such as the book setting and the complex terrain scenario. This is because random placement of fingertips on the terrain can result in a mismatch between the trajectory and the height of the terrain, increasing the likelihood of the finger being stuck or losing contact.

The evaluation results in Fig. 9 show that the PP-Tac pipeline outperforms all baselines. Specifically, comparisons between tactile control and the open-loop trajectory replay method were conducted. The open-loop baseline achieved a lower success rate compared to PP-Tac, highlighting the importance of tactile-based grasping for online force control. We also observed that the teleoperation baseline using a gripper achieved some successful cases in grasping cloth and plastic bags, albeit with lower performance than PP-Tac. This is due to the ease of detecting the initial grasp point on

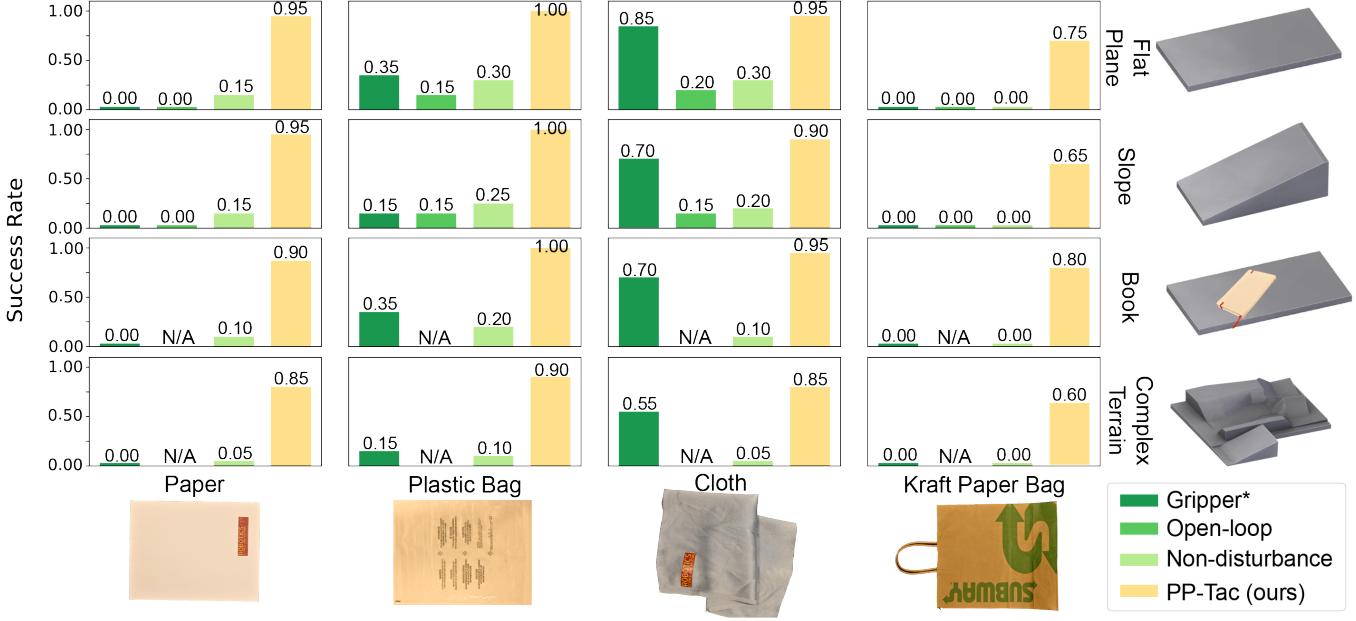


Fig. 9: Experimental Results. Evaluations were conducted to quantify the success rate of grasping four different flat objects (paper, plastic bag, cloth, and paper bag) across four terrain setups (plane, slope, book placement, and randomly generated complex terrain). Baseline conditions included: (1) Gripper*: grasp using a bi-finger gripper controlled by teleoperation; (2) Open-loop: grasp using our dexterous hand in open-loop (without tactile feedback); (3) Non-disturbance: grasp using our dexterous hand with tactile sensors, where the diffusion policy was trained without domain randomization disturbances; and (4) PP-Tac(ours): grasp using our full PP-Tac pipeline. Each condition was repeated 20 times. Note that open-loop grasp control is not feasible on uncertain terrains, and these cases are labeled as ‘N/A’.

these soft materials through human perception, and combined with human intelligence enabling grasp adjustments through visual feedback. However, for stiffer materials like paper and kraft paper, the bi-finger gripper failed completely. Therefore, we conclude that the PP-Tac pipeline is the most suitable configuration for handling flat objects.

E. Ablation Studies

1) *Influence of Material Stiffness:* We found that the material’s stiffness (represented by its thickness), significantly influences the task’s success rate. To demonstrate this effect, we created flat objects by stacking paper pages bonded with adhesive. The experimental results are shown in Tab. I. As the number of paper pages increased, the grasp success rate decreased significantly. Additionally, the increase in material stiffness also led to a higher number of detected slips.

2) *Influence of Data Disturbance:* We emphasize the importance of the data disturbance technique for domain randomization (introduced in Section V-B). To quantify its impact, we conducted ablation studies comparing grasp performance before and after adding four types of disturbances to the prefix motion x^{prefix} . Experimental results demonstrate that this technique significantly enhances performance. As shown in the “Non-disturbance” baseline in Section VI-C, removing data disturbance led to a notable performance drop across all experiments, often resulting in complete failure when grasping stiff objects, such as kraft paper bags. This underscores the improved generalization and higher grasp success rates

enabled by domain randomization. However, a drawback of this technique is the increased training time, requiring approximately 400,000 additional iterations to achieve the same loss as training without data disturbance.

VII. LIMITATIONS

We have observed the following limitations in our system. One limitation is determining the initial force (sensor’s target deformation depth) required for successful grasping. While our algorithm can adaptively adjust the force magnitude online, an appropriate initial value must still be manually set, which remains an empirical parameter-tuning process. If the initial value is too small, the grasp is more likely to fail due to the additional time and finger sliding distance needed for adaptation to a reasonable value. Conversely, if the initial value is too large, excessive friction may exceed the load capacity of the hand motors. In addition to the initial value, the adaptive algorithm for adjusting force also has room for improvement, particularly with highly stiff materials such as kraft paper bags on non-flat surfaces.

VIII. CONCLUSIONS

This paper presents PP-Tac, a coordinated hand-arm system designed to manipulate thin, flat objects such as paper and fabric. The system is equipped with a multi-fingered, vision-based tactile sensor that is easy to fabricate and deploy on the hand’s fingertips. The sensor can detect contact on its curved, omnidirectional surfaces, enabling the system to measure force and friction during contact. This capability helps minimize

slip and increases the likelihood of material deformation when handling flat materials. Based on this hand design, the grasping motion is planned using a data-driven approach. We developed an efficient synthesis algorithm to generate sliding trajectories across various terrain shapes and sensor deformation conditions, resulting in a dataset of 500,000 trajectory samples. Using this dataset and a domain randomization technique, we trained a diffusion policy that enables adaptation to diverse terrains in real-world settings. Experimental results show that our system can successfully grasp flat objects of varying thicknesses and stiffness, achieving a success rate of 87.5%. Additionally, the proposed policy demonstrates robustness to external disturbances and adapts well to different support terrain surfaces.

APPENDIX

A. Detail of Camera Calibration

In this section, we introduce the camera calibration process as part of the overall sensor calibration. Since the tactile sensor is enclosed by an opaque, rounded membrane, conventional calibration board methods cannot be used to determine the pinhole camera's extrinsic parameters. To address this, we designed an indentation setup (as shown in Fig. 10) to capture a sufficient number of spatial points in a known sensor frame, identify their corresponding 2D-pixel coordinates in the image, and establish the mapping between the sensor frame and the image frame. First, the camera's intrinsic parameters K was obtained, either from the camera manufacturer or calibrated using high-precision calibration boards [45]. Next, we define a three-dimensional coordinate system, referred to as the sensor frame (x, y, z) with its origin at the center of the elastomer, as shown in Fig. 10(a). To facilitate the calibration, A custom 3D-printed holder secures the sensor (Fig. 10(b)), while another 3D-printed hemispherical indicator is attached to the holder's groove (Fig. 10(c)). Small pins with a diameter of 1.5mm, serving as indenters, are inserted into pre-defined holes within the indicator for 28 trails. For each trail, the contact positions are recorded both in the camera image as $p_{ij} = (u_{ij}, v_{ij})$ and in the sensor frame as $P_{i,j} = (x_{ij}, y_{ij}, z_{ij})$, where i denotes the trail index and j denotes the contact point index within the trail. The contact positions in the camera image are detected by subtracting the captured image from a reference image without indentation. We use solvePnP [22] to calculate the extrinsic parameters that includes rotation matrix A and translation vector b such that:

$$p_{ij} = K[A \mid b]P_{i,j} \quad (8)$$

After obtaining the intrinsic and extrinsic parameters of the camera, we can project the sensor's curved surface from the sensor frame onto the image frame, obtaining the sensor surface reference projection D (Eq. (9)), by which the depth value on the pixel (u, v) can be queried. This process can be efficiently implemented using the Pyrender [1].

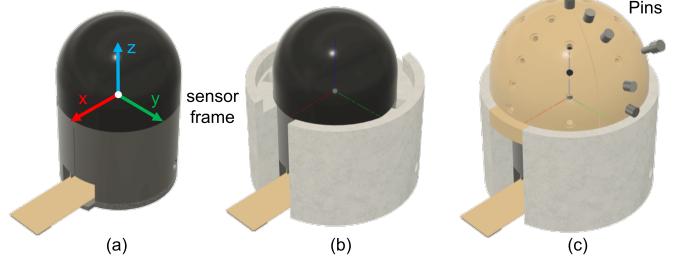


Fig. 10: **Camera calibration using an indentation setup:** The sensor frame is first defined in (a). A holder is designed and 3D-printed to secure the sensor, as shown in (b). A hemispherical indicator is designed and 3D-printed to attach to the sensor holder. Pins are inserted into pre-defined holes to serve as indenters for recording contact locations in the sensor frame, as shown in (c).

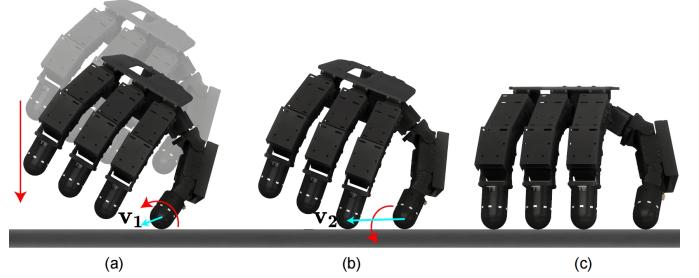


Fig. 11: **Example of establishing contact:** First, the hand descends until a finger makes contact with the surface. A fixed-point rotation is performed around the contacting finger, as shown in (a). The hand then continues to rotate until a second finger makes contact, triggering a fixed-axis rotation around both contacting fingers, as shown in (b). The process is complete when three or more fingers are in contact, as shown in (c).

$$D(u, v) = \begin{bmatrix} u \\ Z_c K^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \end{bmatrix}_{[3,:]}, \quad (9)$$

where $[u \ v \ 1]^T$ and Z_c are given as:

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} \frac{(A[x,y,z]^T + b)_x}{Z_c} \\ \frac{(A[x,y,z]^T + b)_y}{Z_c} \\ 1 \end{bmatrix}, Z_c = (A[x, y, z]^T + b)_z. \quad (10)$$

B. Detail of Establish Contact

In this section, we detail our approach to generate contact with a flat object using the fingertips. The goal is to control the hand to ensure that at least three fingertips are in contact with the surface. We denote the four fingertips as f_t (thumb), f_i (index), f_m (middle), and f_r (ring). The contact states are represented by two sets: C , which includes the fingers in contact, and N , which includes the fingers not in contact. The complete process is illustrated in Fig. 11.

1) **Establish First Contact:** Starting from status when all fingers are hovering (i.e., $C = \{\phi\}$, $N = \{f_t, f_i, f_m, f_r\}$), the hand is controlled to move downward till one finger touches the surface. For example, if the thumb touches the surface (Fig. 11), the contact state sets are updated to $C = \{f_t\}$, $N = \{f_i, f_m, f_r\}$

2) **Establish Second Contact:** Once the first contact is made, the hand rotates around the first finger's contact point to create the second contact point. To achieve this, we first obtain the centroid point of the fingertip in contact (denoted as (x_c, y_c, z_c)), and compute the centroid point of fingertip positions in N (denoted as (x_n, y_n, z_n)). This allows us to calculate the rotational axis as:

$$\mathbf{v}_1 = R_z(90^\circ)(x_n - x_c, y_n - y_c, z_n - z_c)^T, \quad (11)$$

where $R_z(90^\circ)$ is the rotation matrix for a 90-degree rotation around the z-axis. Given θ , \mathbf{v}_1 calculated before, robot arm's target end-effector pose ${}^b_{ee'}T$ leading to such rotation can be obtained via Rodrigues' rotation formula:

$$R(\theta, \mathbf{v}_1) = I + \sin(\theta) \begin{bmatrix} 0 & -v_{1z} & v_{1y} \\ v_{1z} & 0 & -v_{1x} \\ -v_{1y} & v_{1x} & 0 \end{bmatrix} + (1 - \cos(\theta)) \begin{bmatrix} 0 & -v_{1z} & v_{1y} \\ v_{1z} & 0 & -v_{1x} \\ -v_{1y} & v_{1x} & 0 \end{bmatrix}^2, \quad (12)$$

The target end effector pose of the robot arm can be calculated as:

$${}^b_{ee'}T = {}^b_{ee'}T {}^{ee'}_c T {}^c_{c'} \hat{T} {}^{c'}_{ee'}T, \quad (13)$$

$${}^{c'}_{c''} \hat{T} = \begin{bmatrix} R(\theta, \mathbf{v}_1) & 0 \\ 0 & 1 \end{bmatrix}, \quad (14)$$

where b denotes the base of the robot arm, ee and ee' represent the end effector before and after the movement, and c and c' represent the positions (x_c, y_c, z_c) before and after the rotation. The robot arm is then controlled to gradually increase θ until the second fingertip contacts the object surface. Once this occurs, we update the contact states to $C = \{f_t, f_i\}$ and $N = \{f_m, f_r\}$.

3) **Establishing Third Contact:** In this step, the hand rotates around an axis defined by the first and second contact points until the third fingertip makes contact. For instance, if the thumb and index finger make contact, the rotation axis is $\mathbf{v}_2 = \overrightarrow{f_t f_i}$. The arm's target end-effector pose for this rotation is:

$${}^{ee''}_{ee'}T = {}^{ee'}_{ee'}T {}^{ee'}_{c''} T {}^{c''}_{c'} \hat{T} {}^{c'}_{ee''}T, \quad (15)$$

$${}^{c''}_{c''} \hat{T} = \begin{bmatrix} R(\theta', \mathbf{v}_2) & 0 \\ 0 & 1 \end{bmatrix}, \quad (16)$$

where c'' and ee'' are c' and ee' after rotation specified by \mathbf{v}_2 . During execution, the angle θ' is gradually increased until a new fingertip contacts the surface, achieving the desired target end-effector pose ${}^b_{ee'}T$. Note that these steps may not always be required. In some cases, we observe that the third finger may already be in the contact state when we attempt to establish contact with the second finger.

C. List of Symbols

The definition of symbols can be found in Tab. II.

TABLE II: Summary of symbols and notations

Symbols	Descriptions
u, v	Pixel coordinates in VBTS.
X_c, Y_c, Z_c	Camera coordinates in VBTS.
x, y, z	Sensor coordinates in VBTS.
K	The intrinsic parameters of the camera in VBTS.
A, b	The extrinsic parameters of the camera in VBTS.
D	Sensor surface reference projection in VBTS.
M	Depth mapping function in VBTS.
q	Rotation angle of controllable hand joints.
\dot{q}	Angular velocity of controllable hand joints.
p_j	Positional coordinate of hand joints in arm's base axis.
\dot{p}_j	Linear velocity of hand joints in arm's base axis.
R	Wrist's (end effector of arm) 6D rotation.
Ω	Angular velocity of hand pose.
p_{wrist}	Wrist (end-effector of arm)'s height along arm's z-axis.
\dot{p}_{wrist}	Linear velocity of $p_{ee'}$.
d_{tac}	The deformation depth readings from four fingertip tactile sensors.
\bar{d}_{tac}	The target deformation depth.
\mathcal{D}	State variable's dimension.
γ	Hand joint angles $q^{1:N_{data}}$, wrist's (end effector of arm) 6D rotation $R^{1:N_{data}}$ and wrist's translation along z-axis $p_{ee'}^{1:N_{data}}$ for overall trajectory.
N_{data}	Length of synthesis motion sequence.
N_{pred}	Length of predicted actions.
x_{pred}	Future motion predicted by PP-Tac policy.
N_{prefix}	Length of historical actions.
x_{prefix}	The historical action sequence.
t	Diffusion step.

REFERENCES

- [1] Pyrender. <https://github.com/mmatl/pyrender>, 2020.
- [2] Enrique Amigó, Julio Gonzalo, and Felisa Verdejo. A general evaluation measure for document organization tasks. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*, pages 643–652, 2013.
- [3] Veronica E Arriola-Rios and Jeremy L Wyatt. A multimodal model of object deformation under robotic pushing. *IEEE Transactions on Cognitive and Developmental Systems*, 9(2):153–169, 2017.
- [4] Aude Billard and Danica Kragic. Trends and challenges in robot manipulation. *Science*, 364(6446):eaat8414, 2019.
- [5] Tara Boroushaki, Junshan Leng, Ian Clester, Alberto Rodriguez, and Fadel Adib. Robotic grasping of fully-occluded objects using rf perception. In *2021 IEEE*

- International Conference on Robotics and Automation (ICRA)*, pages 923–929. IEEE, 2021.
- [6] Arkadeep Narayan Chaudhury, Timothy Man, Wenzhen Yuan, and Christopher G. Atkeson. Using collocated vision and tactile sensors for visual servoing and localization. *IEEE Robotics and Automation Letters*, 7(2):3427–3434, 2022. doi: 10.1109/LRA.2022.3146565.
 - [7] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023.
 - [8] Xinke Deng, Yu Xiang, Arsalan Mousavian, Clemens Eppner, Timothy Bretl, and Dieter Fox. Self-supervised 6d object pose estimation for robot manipulation. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3665–3671. IEEE, 2020.
 - [9] Won Kyung Do and Monroe Kennedy. Densetact: Optical tactile sensor for dense shape reconstruction. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 6188–6194. IEEE, 2022.
 - [10] Won Kyung Do, Bianca Aumann, Camille Chungyoun, and Monroe Kennedy. Inter-finger small object manipulation with densetact optical tactile sensor. *IEEE Robotics and Automation Letters*, 2023.
 - [11] Mehmet Remzi Dogar and Siddhartha S Srinivasa. A framework for push-grasping in clutter. In *Robotics: Science and systems*, volume 2, 2011.
 - [12] Siyuan Dong, Daolin Ma, Elliott Donlon, and Alberto Rodriguez. Maintaining grasps within slipping bounds by monitoring incipient slip. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 3818–3824. IEEE, 2019.
 - [13] Satoshi Funabashi, Tomoki Isobe, Shun Ogasa, Tetsuya Ogata, Alexander Schmitz, Tito Pradhono Tomo, and Shigeki Sugano. Stable in-grasp manipulation with a low-cost robot hand by using 3-axis tactile sensors with a cnn. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 9166–9173. IEEE, 2020.
 - [14] Rafael Herguedas, Gonzalo López-Nicolás, Rosario Aragüés, and Carlos Sagüés. Survey on multi-robot manipulation of deformable objects. In *2019 24th IEEE International Conference on Emerging Technologies and Factory Automation (ETFA)*, pages 977–984. IEEE, 2019.
 - [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf.
 - [16] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models, 2022.
 - [17] Mohsen Kaboli, Rich Walker, Gordon Cheng, et al. In-hand object recognition via texture properties with robotic hands, artificial skin, and novel tactile descriptors. In *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, pages 1155–1160. IEEE, 2015.
 - [18] Justin Kerr, Huang Huang, Albert Wilcox, Ryan Hoque, Jeffrey Ichnowski, Roberto Calandra, and Ken Goldberg. Self-supervised visuo-tactile pretraining to locate and follow garment features. In *Proceedings of Robotics: Science and Systems*, 2023.
 - [19] 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. In *Proceedings of Robotics: Science and Systems*, 2023.
 - [20] Michael Krawez, Tim Caselitz, Jugesh Sundram, Mark Van Loock, and Wolfram Burgard. Real-time outdoor illumination estimation for camera tracking in indoor environments. *IEEE Robotics and Automation Letters*, 6(3):6084–6091, 2021.
 - [21] Mike Lambeta, Tingfan Wu, Ali Sengul, Victoria Rose Most, Nolan Black, Kevin Sawyer, Romeo Mercado, Haozhi Qi, Alexander Sohn, Byron Taylor, et al. Digitizing touch with an artificial multimodal fingertip. *arXiv preprint arXiv:2411.02479*, 2024.
 - [22] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. Ep n p: An accurate o (n) solution to the p n p problem. *International journal of computer vision*, 81:155–166, 2009.
 - [23] Mengdi Li, Cornelius Weber, Matthias Kerzel, Jae Hee Lee, Zheni Zeng, Zhiyuan Liu, and Stefan Wermter. Robotic occlusion reasoning for efficient object existence prediction. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2686–2692. IEEE, 2021.
 - [24] Rui Li, Robert Platt, Wenzhen Yuan, Andreas ten Pas, Nathan Roscup, Mandayam A. Srinivasan, and Edward Adelson. Localization and manipulation of small parts using gelsight tactile sensing. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3988–3993, 2014. doi: 10.1109/IROS.2014.6943123.
 - [25] Yinxiao Li, Danfei Xu, Yonghao Yue, Yan Wang, Shih-Fu Chang, Eitan Grinspun, and Peter K Allen. Re-grasping and unfolding of garments using predictive thin shell modeling. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1382–1388. IEEE, 2015.
 - [26] Yunzhu Li, Jiajun Wu, Russ Tedrake, Joshua B Tenenbaum, and Antonio Torralba. Learning particle dynamics for manipulating rigid bodies, deformable objects, and fluids. *arXiv preprint arXiv:1810.01566*, 2018.
 - [27] Changyi Lin, Ziqi Lin, Shaoxiong Wang, and Huazhe

- Xu. Dtact: A vision-based tactile sensor that measures high-resolution 3d geometry directly from darkness. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 10359–10366. IEEE, 2023.
- [28] Changyi Lin, Han Zhang, Jikai Xu, Lei Wu, and Huazhe Xu. 9dtact: A compact vision-based tactile sensor for accurate 3d shape reconstruction and generalizable 6d force estimation. *IEEE Robotics and Automation Letters*, 2023.
- [29] Xiaofei Liu, Wuqiang Yang, Fan Meng, and Tengchen Sun. Material recognition using robotic hand with capacitive tactile sensor array and machine learning. *IEEE Transactions on Instrumentation and Measurement*, 2024.
- [30] Ashvin Nair, Dian Chen, Pulkit Agrawal, Phillip Isola, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Combining self-supervised learning and imitation for vision-based rope manipulation. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 2146–2153. IEEE, 2017.
- [31] Kei Ota, Devesh K Jha, Hsiao-Yu Tung, and Joshua Tenenbaum. Tactile-Filter: Interactive Tactile Perception for Part Mating. In *Proceedings of Robotics: Science and Systems*, Daegu, Republic of Korea, July 2023. doi: 10.15607/RSS.2023.XIX.079.
- [32] 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.
- [33] Wonik Robotics. Allegro Hand, 2024. URL <https://www.allegrohand.com/ah-v4-main>.
- [34] Brian Scassellati, Henny Admoni, and Maja Matarić. Robots for use in autism research. *Annual review of biomedical engineering*, 14(1):275–294, 2012.
- [35] Yu She, Shaoxiong Wang, Siyuan Dong, Neha Sunil, Alberto Rodriguez, and Edward Adelson. Cable manipulation with a tactile-reactive gripper. *The International Journal of Robotics Research*, 40(12-14):1385–1401, 2021.
- [36] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [37] Li Sun, Gerardo Aragon-Camarasa, Simon Rogers, and J Paul Siebert. Accurate garment surface analysis using an active stereo robot head with application to dual-arm flattening. In *2015 IEEE international conference on robotics and automation (ICRA)*, pages 185–192. IEEE, 2015.
- [38] Sudharshan Suresh, Haozhi Qi, Tingfan Wu, Taosha Fan, Luis Pineda, Mike Lambeta, Jitendra Malik, Mrinal Kalakrishnan, Roberto Calandra, Michael Kaess, et al. Neuralfeels with neural fields: Visuotactile perception for in-hand manipulation. *Science Robotics*, 9(96):eadl0628, 2024.
- [39] Ian H Taylor, Siyuan Dong, and Alberto Rodriguez. Gelslim 3.0: High-resolution measurement of shape, force and slip in a compact tactile-sensing finger. In *International Conference on Robotics and Automation (ICRA)*, pages 10781–10787, 2022.
- [40] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=SJ1kSyO2jwu>.
- [41] Megha H Tippur and Edward H Adelson. Rainbowsight: A family of generalizable, curved, camera-based tactile sensors for shape reconstruction. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1114–1120. IEEE, 2024.
- [42] Benjamin Ward-Cherrier, Nicholas Pestell, Luke Cramphorn, Benjamin Winstone, Maria Elena Giannaccini, Jonathan Rossiter, and Nathan F Lepora. The tactip family: Soft optical tactile sensors with 3d-printed biomimetic morphologies. *Soft robotics*, 5(2):216–227, 2018.
- [43] Xiaolong Wu and Cédric Pradalier. Illumination robust monocular direct visual odometry for outdoor environment mapping. In *2019 International conference on robotics and automation (ICRA)*, pages 2392–2398. IEEE, 2019.
- [44] Mengyuan Yan, Yilin Zhu, Ning Jin, and Jeannette Bohg. Self-supervised learning of state estimation for manipulating deformable linear objects. *IEEE robotics and automation letters*, 5(2):2372–2379, 2020.
- [45] Zhengyou Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2002.
- [46] Yi Zhou, Connelly Barnes, Lu Jingwan, Yang Jimei, and Li Hao. On the continuity of rotation representations in neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [47] Jihong Zhu, Andrea Cherubini, Claire Dune, David Navarro-Alarcon, Farshid Alambeigi, Dmitry Berenson, Fanny Ficuciello, Kensuke Harada, Jens Kober, Xiang Li, et al. Challenges and outlook in robotic manipulation of deformable objects. *IEEE Robotics & Automation Magazine*, 29(3):67–77, 2022.