

Iterative Interactive Modeling for Knotting Plastic Bags

Anonymous Author(s)

Affiliation

Address

email

Abstract:

Deformable object manipulation has great research significance for the robotic community and numerous applications in daily life. In this work, we study how to knot plastic bags that are randomly dropped from the air with a dual arm robot based on image input. The complex initial configuration and intricate material and dynamics properties of plastic bags pose challenges for reliable perception and planning. Directly knotting it from random initial states is difficult. To tackle this problem, we propose Iterative Interactive Modeling (IIM) to first adjust the plastic bag to a standing pose with imitation learning to establish a high-confidence key-point skeleton model, then perform a set of learned motion primitives to knot it. We leverage spatial action maps to accomplish the iterative pick-and-place action and a graph convolutional network to evaluate the adjusted pose during the IIM process. In experiments, we achieve an 85.0% success rate in knotting 4 different plastic bags including one that has no demonstration.

Keywords: Plastic Bag Manipulation, Learning from Demonstrations

1 Introduction

Deformable object manipulation (DOM) has been a long standing problem in robotics. Researchers have been studying manipulating various kinds of deformable objects, from linear objects [1, 2], fabrics [3, 4], papers [5] to elastic and elasto-plastic objects [6, 7, 8]. Apart from them, plastic bag, perhaps the most widely used application of plastics in our daily life, has remained unexplored in robotics literature due to the high complexity involved in modeling and controlling its deformation. Endowing robots with the ability to manipulate plastic bags can spawn diverse industrial and domestic applications in warehouses, garbage dumps, and supermarkets. Knotting plastic bags is one of the most representative manipulation tasks on plastic bags for a robot to show its dexterity in the real world. In this paper, we study how to knot plastic bags that are randomly dropped from the air with a dual arm robot based on raw image input.

A randomly dropped plastic bag is shown in Figure 1(a). Knotting it from such an irregular initial configuration requires locating and taking the handles out from the mess and delicate coordination of dual arms to tie the knot. Compared to previous knot tying tasks that have been

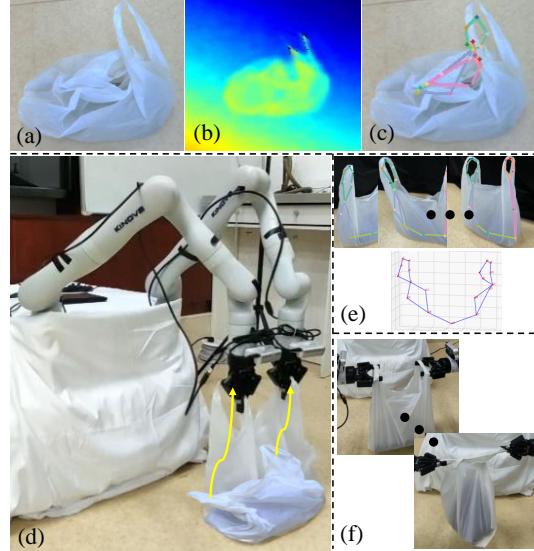


Figure 1: Overview of our work. (a) A randomly dropped plastic bag. (b) Inaccurate depth measurement from Intel Realsense D435i. (c) Unstable keypoint detection at the initial state. (d)(e) We propose IIM to first iteratively interact with the plastic bag to shape it to a standing pose, then build a high-confidence 3D keypoint skeleton with triangulation. (f) Knot tying with a set of learned motion primitives.

well studied on ropes [1, 9, 10, 11], this task poses

44 new challenges for both perception and planning, because of the complex initial configuration and
45 the intricate material and dynamics property of plastic bags. For perception, as shown in Figure
46 1(b), the translucent and non-Lambertian surface of plastic bags makes depth sensing inaccurate,
47 thus 3D vision based models such as point cloud [12, 13, 14] or voxels [15] are not suitable for
48 our task. Meanwhile, 2D vision models such as 2D keypoints [16] are not reliable due to the se-
49 vere self-occlusion and partial observation problem caused by complex initial configurations (Figure
50 1(c)). For planning, mathematical physical models are difficult to build for plastic bags since their
51 physical parameters are very hard to estimate, and to the best of our knowledge, no simulator can
52 be employed. Learning a dynamics model for planning [7, 13, 15] is also infeasible since motions
53 of plastic bags are not quasi-static. Although some image-based imitation or reinforcement learning
54 approaches [3, 4, 17, 18] can directly learn a visual policy without a visual representation or physical
55 model, their learned policies either employ just a small action space such as pick-and-place that is
56 not enough to accomplish the knot tying task, or are unable to cope with complex initial configu-
57 rations that appear in a randomly dropped plastic bag. Therefore, directly knotting the plastic bag
58 from initial configurations is hard for both perception and planning.

59 Fortunately, the elasto-plastic property of plastic bags gives us another way to accomplish the knot
60 tying task. Elasto-plastic objects can be shaped to a predefined target state in equilibrium and un-
61 changing. This is called deformation control [19], and previous works have studied how to learn a
62 dynamics model [13] or estimate physical parameters [20] of elasto-plastic objects such as sponges,
63 plasticine, or clay by interacting with them during this shaping process. However, for our task, this
64 interactive shaping process is more valuable for reducing difficulties of perception and planning:
65 shaping a plastic bag from a randomly dropped state to a standing pose (Figure 1(d)) is a process of
66 straightening it, and the partial observation and self-occlusion problem can be gradually alleviated
67 during this process. This allows us to build a more reliable visual model (such as keypoint skeleton,
68 Figure 1(e)), and the knot tying task can be realized from the standing pose based on this model.
69 More importantly, the plastic bag does not need to be exactly shaped to a standard registered pose:
70 we only need it to be spread enough to support building a reliable visual model. This allows the
71 shaping process itself to avoid complex perception or planning, thus it can be accomplished by com-
72 plete learning approaches with only iterative pick-and-place actions [17, 4, 21]. Therefore, we aim
73 to develop a visual learning policy to iteratively adjust the plastic bag from randomly dropped initial
74 configurations, and an evaluation module to evaluate if the shaped pose is good enough for reliable
75 perception and the down-streaming knot tying task.

76 In this work, we propose Iterative Interactive Modeling (IIM) for knotting plastic bags randomly
77 dropped from the air with only image input. We train the robot to first shape the plastic bag to
78 a standing pose with the help of demonstrations and then establish a keypoint skeleton model with
79 multi-view stereo images to knot it with a set of learned adaptive motion primitives. Specifically, the
80 robot iteratively performs different kinds of top-down pick-and-place actions on part of the plastic
81 bag to outspread it, meanwhile, a task progress module evaluates if the pose is good enough for knot
82 tying. We leverage spatial action maps [17, 18] to accomplish the pick-and-place action, and train a
83 graph convolutional network as the task progress module with the same demonstrations to evaluate
84 the keypoint skeleton during the adjusting process. To enable keypoint detecting on plastic bags,
85 we provide the first 2D plastic bag keypoint dataset *PBPose* with 43,200 images to train an off-the-
86 shelf 2D keypoint detection model RLE [22]. After the IIM process, we lift the plastic bag into air
87 with geometrically constrained planning to knot the plastic bag with a set of motion primitives with
88 trained action parameters from CNN. In our experiments, we achieve an 85.0% success rate in tying
89 four different plastic bags (one of them has no demonstration) that randomly dropped from the air
90 with our dual Kinova Gen3 arms equipped with standard Robotiq 2F-85 grippers, with only 100
91 demonstrations (1.5 hours) provided for each plastic bag. In summary:

- 92 • We propose Iterative Interactive Modeling (IIM) for complex elasto-plastic object manipu-
93 lation that iteratively shapes the object to facilitate more reliable perception and planning.
- 94 • We leverage spatial action maps, graph convolutional networks and the RLE model to per-
95 form IIM on plastic bags, and train motion primitives to accomplish the knot tying task.
- 96 • We build the first dual-arm robotic system to knot plastic bags randomly dropped from the
97 air with the provided *PBPose* dataset and a small number of demonstrations.

98 2 Related Works

99 Deformable object manipulation (DOM) has long been a challenging area of robotics research. The
100 challenges of DOM come from two properties of deformable objects: the infinite degrees of freedom

101 and the complex non-linear dynamics, which lead to difficulties for perception, modeling, and planning.
102 In this section, we compare different techniques for modeling and manipulating deformable
103 objects and discuss whether they are suitable for knotting plastic bags.

104 2.1 Visual Representations of Deformable Objects

105 Keypoints [23, 24, 25] are commonly used representations. They are sparse representations for the
106 structure of target objects, and grasping points are usually generated from them. Keypoint detec-
107 tors can be trained either with supervised or unsupervised ways [26] to find keypoints with image
108 or point cloud inputs to allow robots to manipulate the object based on some geometry loss on the
109 keypoints. Dense visual descriptors [8, 27, 28, 29] are recently used as a real-time pixel-wise dense
110 representation for many kinds of objects, and humans can specify grasping points on it. Besides, the
111 3D vision community has made great progress in modeling deformable objects, with voxels [30],
112 meshes [31], convexes [32], or implicit functions such as flow-based model [33] or neural radiance
113 fields [34], and most of these model can represent the deformable object in high fidelity. Other
114 works seek to use particles [35, 14, 36] that are down-sampled from point cloud data to represent
115 deformable objects. These methods usually combine graph neural networks and differentiable sim-
116 ulators [37, 38] to learn the dynamics of objects. Lastly, visual policy learning methods simply use
117 high-dimensional feature embedding learned from deep neural networks as object representations
118 for manipulation [3, 4, 11, 39] and have achieved great success in various tasks.

119 For plastic bags, current commercial depth sensors fail to give accurate depth estimation because of
120 the translucent and non-Lambertian surface, thus 3D vision based methods are not feasible. Since
121 we aim to manipulate the plastic bag rather than studying how to build a refined representation,
122 those high-fidelity representations are not necessary: they all degenerate to grasping points [28] in
123 practice for manipulation. Thus, we choose keypoint skeleton as the representation of plastic bags.

124 2.2 Manipulating Deformable Objects

125 A lot of work aims to build a dynamics model for the target object first and then perform motion
126 planning on it to manipulate the object. Conventional works build the dynamics model mathe-
127 matically [40, 41, 42], and some of them have achieved promising results recently by planning
128 on simpler approximated dynamics and using local controllers to tackle actual complex dynamics
129 [43, 44, 45]. Other works aim to learn the dynamics model from visual data [13, 15, 7] or tactile
130 data [46, 47]. Recently, reinforcement learning and imitation learning methods are developed to di-
131 rectly manipulate the object without a dynamics model, and have achieved success on various tasks
132 [3, 48, 1, 10, 17, 49]. Some of them train the robot in simulation with fast virtual experiences and
133 then transfer the learned policy directly to real robots or with sim-to-real methods [27, 50]. Others
134 try to directly train the robot in the real world with its own experience or expert demonstrations [4].

135 For manipulating plastic bags, dynamics models are hard to build since the complex non-linear
136 dynamics, and the motion of plastic bag is not always quasi-static: some parts of the plastic bag may
137 collapse after it has been moved. In this work, we follow the idea of most visual policy learning
138 approaches that directly learn a pick-and-place policy with a limited action space to accomplish the
139 adjusting task, and use geometrically constrained planning to lift the plastic bag to the air to knot it
140 with a set of learned motion primitives.

141 3 Method

142 The goal of this paper is to enable a dual arm robotic system to knot plastic bags from initial ran-
143 domly dropped states. We use keypoint skeleton as the visual representation for plastic bags. The
144 robot is trained to iteratively adjust the plastic bag to a standing pose to build a complete and high-
145 confidence keypoint model with the help of a task progress module, then tie the knot with a set of
146 learned motion primitives. We introduce the problem setting in 3.1, the keypoint detection model in
147 3.2, the iterative interactive modeling process in 3.3, and the knotting process in 3.4.

148 3.1 Problem Statement

149 We use the most common type of plastic bags that are colored, medium-sized (33.7~42.5 cm wide
150 and 40.2~42.3 cm long when laid out), translucent, and have two looped handles, as shown in Figure
151 2(a). We choose four kinds of plastic bags that vary in size and color. In order to make the tying
152 problem possible, we put some items in the plastic bag to resist indoor airflow. The plastic bag is
153 randomly dropped from the air to form an initial state. The knot we tie in this work is a kind of *Ian
Knot* [9] which can be tied with two motion primitives (Action 1 and Action 2) from the stage of

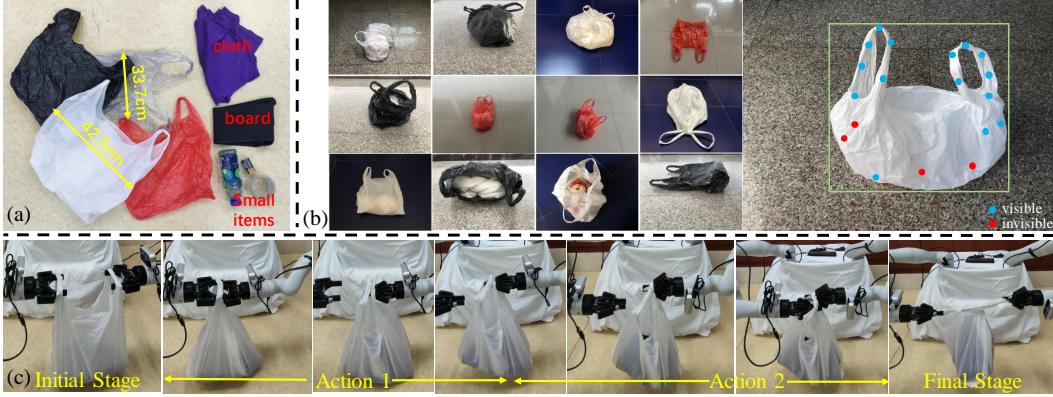


Figure 2: (a) Plastic bags and fillers used in this work. We give demonstrations on white, red and black bags. The grey plastic bag is used to show the generalization of IIM. We fill plastic bags with different items varying in size and shape, including cloths, paperboards, bottles, tape volumes, and dry batteries. (b) The *PBPose* dataset. Images are captured in various backgrounds, knotting states, occupied states, illumination conditions, placement status (standing or lying down), and angles of view. (c) We design two motion primitives for tying an *Ian Knot* in the air.

155 hanging on grippers, as illustrated in Figure 2(c): 1) The right arm grasps the rear side of the left
 156 handle to allow the left arm to go back; 2) The right arm rotates 90° to let the left arm grasp the front
 157 side of the right handle. The knot is tightened by pulling two handles away from each other.

158 The whole experimental setup is shown in Figure 3. We use two identical Kinova Gen3 (6DoF)
 159 arms and Robotiq 2F-85 grippers for our experiments. We mount two Intel Realsense D435i cameras
 160 C_{left} and C_{right} on the end of each arm respectively to get image inputs I_{left} and I_{right} from
 161 the end-effector perspective. The distance between the bases of two arms is 51.2cm. C_{right} looks
 162 straight down at the XY plane from a top view to get I_{right} within the range of 75.4cm × 56.6cm
 163 with the plastic bag at the center. On the other side, we get I_{left} by moving C_{left} aiming at the
 164 plastic from a set of widely-separated poses to get stereo images.

165 3.2 3D Keypoint Skeleton for Plastic Bags

166 Directly 3D keypoint detecting methods and reconstruction methods are not suitable for plastic bags
 167 as discussed in section 2.1. Thus in this work we detect 2D keypoints of plastic bags with image
 168 input and reconstruct 3D keypoint skeleton with multi-view geometry. This requires a plastic bag
 169 keypoint dataset and a well-trained 2D keypoint detection model.

170 To this end, we provide *PBPose*, the first 2D plastic bag keypoint dataset with 43,200 images of four
 171 different kinds of plastic bags. Images in *PBPose* vary from the background, knotting state, occupied
 172 state, illumination condition, placement status, and angle of view, as shown in Figure 2(b). For each
 173 image, we manually label 19 keypoints of the plastic bag from handle to bottom to establish a full
 174 skeleton of it, along with point visible properties and the bounding box. Based on this dataset, we
 175 train a RLE network [22] f for 2D keypoint detection. f takes as input of a single frame from I_{left}
 176 and outputs predicted 2D joints of the target plastic bag along with confidences of each point.

177 From the illustrated performances on plastic bags at initial states and standing poses in Figure 1(c)
 178 and 1(e), we can see that it is not robust and reliable to detect keypoints from random initial config-
 179 urations. Although this problem can be mitigated by providing more training data at different initial
 180 configurations, it is hard to cover enough cases of the infinite continuous initial states. Thus in this
 181 work we sidestep this problem by shaping the plastic bag to a standing pose.

182 3.3 Iterative Interactive Modeling

183 We adjust the plastic bag from its initial state by iteratively picking and placing it with expert demon-
 184 strations, and automatically evaluate the pose with a task progress module. We introduce these
 185 modules one by one in the following part. Note we do not require the whole adjusting process to be
 186 quasi-static, which means the plastic bag can partially collapse before it can steadily stand.

187 **Observation and Action Space Definition:** At every step, we get h multi-view images $I_{left}^{mv} =$
 188 $\{I_{left}^0, \dots, I_{left}^{h-1}\}$ from C_{left} ($h = 15$ in this work), and one top-view image I_{right} from C_{right} .

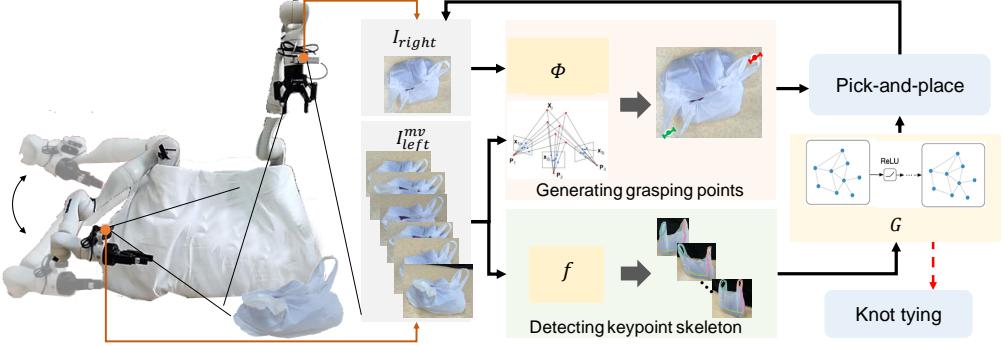


Figure 3: The IIM iteration for knotting plastic bags. At each step, the right camera gets a top-down view image I_{right} , and the left camera gets a sequence of side view images I_{left}^{mv} . We use a spatial action maps module ϕ to get the grasping points, and use a sparse reconstructed point cloud to get the grasping depth. At the same time, we detect 2D keypoints from I_{left}^{mv} with f to facilitate the task progress module G to determine which picking action to use, and when the adjusting pose is good enough to support building a 3D keypoint skeleton to tie the knot.

- 189 The robot can choose two kinds of pick-and-place action at each step: 1) a single-arm picking action
 190 p_1 that goes top-down to location $R \in \mathbb{R}^3$ and picks up to 0.55m height along the Z axis; 2)
 191 a dual-arm picking action p_2 that goes top-down to locations $L, R \in \mathbb{R}^3$ and picks up to 0.55m height
 192 along the Z axis respectively. After each picking, the arm can choose to go to a left/staying/right
 193 placing position to release the plastic bag, as shown in Figure 4. p_1 is used for most pick-and-place
 194 actions and p_2 is only used for final steps to get a good standing pose.
 195 The robot needs to choose which action to use, and grasping and placing locations at every step.
 196 Specifically, the xy locations of grasping points are determined by our learned visual policy. The
 197 z parameter (grasping depth) is determined by first running sparse reconstruction of stereo images
 198 by structure from motion algorithm (SfM) to get sparse point cloud (using COLMAP [51] in this
 199 work), then calculating the average height of the top-10 points on a cylinder with a radius of 1cm
 200 with the center of the bottom surface of the grasping point to get the grasping depth. For determining
 201 the placing location, we segment I_{right} and calculate the relative sizes of the left-side and right-side
 202 areas. If the ratio of two areas is less than 0.5 or greater than 2, the robot will place the plastic bag
 203 to the smaller-side placing point. Otherwise it will release the gripper in situ, as shown in Figure 4.

204 **Demonstration Collection:** A human expert demonstrates how to accomplish the adjusting process
 205 for the robot. For each demonstration trajectory $D = \{p_{c_0}, p_{c_1}, \dots, p_{c_{k-1}}\}$, it includes k picking
 206 actions, where c_0, c_1, \dots, c_{k-1} denote picking categories. For each p_{c_i} , we record I_{left}^{mv} and I_{right}
 207 before performing the action, and the grasping locations $L, R \in \mathbb{R}^2$ from the perspective of C_{right} ,
 208 along with the grasping angles θ_L, θ_R . Thus, $p_{c_i} = \{I_{left}^{mv}, I_{right}, L, R, \theta_L, \theta_R\}_{c_i}$. The expert al-
 209 ways grasps *handles* of the plastic bag for all picking actions. This is important to establish a unified
 210 grasping principle to avoid ambiguity for imitation learning. p_2 may be used several times at the end
 211 of a trajectory to form a good standing pose. In this paper, we find 100 demonstration trajectories
 212 are enough for each kind of plastic bag, which only take 1.5 hour for human to demonstrate.

213 **Visual Grasping Module:** The robot determines the grasping points $\{(L, \theta_L), (R, \theta_R)\}$ from I_{right}
 214 with spatial action maps [17, 21, 18], which have shown promising results for learning heatmaps
 215 of visual-affordances over pixels with fully convolutional networks. They can recover pixel-wise
 216 grasping affordances of different graspings with input images being rotated to achieve a discrete set
 217 of possible actions. Concretely, as shown in Figure 4, given an image from I_{right} , we first segment
 218 out the plastic bag from the background and generate 16 rotated images (multiples of 22.5°), then
 219 pass them through the fully-convolutional network ϕ to predict the corresponding set of heatmaps
 220 within the same size of the input image. In this work we use a pixel resolution of 160×120 . The
 221 pixel with a higher predicted probability among all 16 maps is more suitable for grasping.

222 The problem in our case is that we have to determine one grasping point for p_1 but two grasping
 223 points for p_2 . To this end, we propose a *Selection Rule* based on safe distance constraints to avoid
 224 collisions of grippers: for p_1 and p_2 , the pixel with the highest value corresponds to (R, θ_R) . When
 225 choosing (L, θ_L) for p_2 , we first delete nearby areas (a circle of 5 pixels radius) of (R, θ_R) on all 16

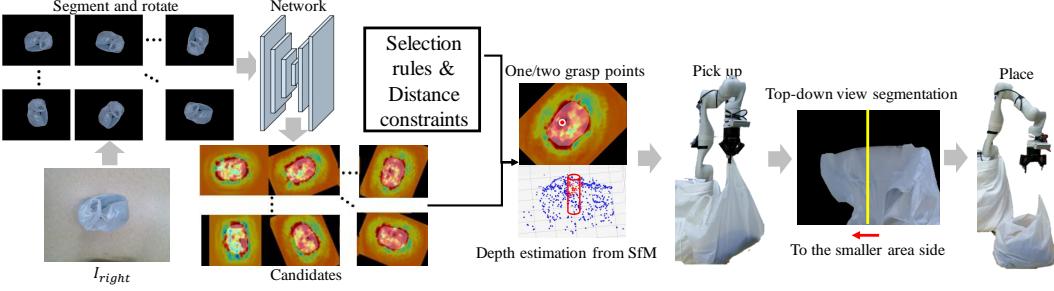


Figure 4: A pick-and-place process. I_{right} is segmented and rotated to get 16 copies. We send these images through network ϕ and choose one or two grasping point(s) from the output candidates according to our selection rules and distance constraints. We get the grasping depth of these points by calculating the depths of top-10 points above the grasping points and perform the picking action. After picking it upward, we determine where to place the plastic bag by segmenting the top-down view and compare relative sizes of left-side and right-side area of the plastic bag.

226 output images, and select the pixel with the highest value in the remaining area. If this pixel has a
227 value higher than 0.8, then we choose it as (L, θ_L) , otherwise we abort p_2 and only perform p_1 .

228 **Task Progress Module:** To evaluate the quality of poses during the adjusting process to determine
229 when the robot can use p_2 and when the pose is good enough to facilitate knot tying, we propose
230 to use a graph convolutional network to classify the keypoint skeleton during the IIM process to
231 predict the *task progress* stages. We represent the 2D keypoint skeleton as a graph formed by state
232 $S = (\mathcal{O}, \mathcal{E})$, where vertices \mathcal{O} are the keypoints. Concretely, $\mathbf{o}_i = \langle \mathbf{x}_i, \mathbf{c}_i \rangle$, where \mathbf{x}_i and \mathbf{c}_i are
233 2D pixel locations and confidences of each keypoint respectively. The edges \mathcal{E} are the connectivity
234 of each pair of keypoints that are predefined. At each step, a graph convolutional network G gets
235 h graphs abstracted from stereo image inputs of I_{left}^{mv} and detected keypoint skeletons by f and
236 calculates the mean state \bar{S} . Then it classifies current step to three categories: 1) ordinary picking
237 step (using p_1); 2) final picking step (using p_2); 3) ending step (no picking action will be chosen).
238 We train G with the standard cross entropy loss and the same demonstrations collected above.

239 After IIM, we get a well-standing plastic bag. We recover the 3D keypoint skeleton of it by performing
240 triangulation from detected 2D keypoints of stereo images of I_{left}^{mv} with the help of calibrated
241 camera intrinsic and extrinsic parameters. The next step is to knot the plastic bag with this model.

242 3.4 In-air Knotting Plastic Bags by Learning Motion Primitives

243 As shown in Figure 2(c), tying an Ian Knot
244 needs to first lift the plastic bag to the air to
245 eliminate the influence from ground. This re-
246quires the direction that allows each arm to go
247 into each ring on handles, as shown in Figure
248 5(a). There are six keypoints on each ring, and
249 they are not on the same plane in \mathbb{R}^3 . Thus
250 we here calculate a *ring plane* that minimizes
251 the total least squares (TLS) distance from each
252 point to it. Then we calculate the perpendicular
253 direction of each ring plane, and let this perpen-
254 dicular go through the center point of the ring,
255 which is the average of 3D coordinates of all the
256 ring points. By TLS, we know this center point
257 is on the ring plane too. We set each robot arm
258 on its perpendicular respectively and place each gripper 0.1m away from the center point of each
259 ring. Then we close both grippers and move the arm along the perpendicular to insert into the ring,
260 and keep both end effectors horizontal with a 0.25m distance to lift the bag to a height of 0.5m.
261 Finally we open both grippers to reach the *initial stage* for tying an Ian Knot.

262 The second step is to perform the two motion primitives. However, as shown in Figure 5(b), the
263 dangling handles of plastic bags on the gripper may hang in different specific places. This requires
264 both actions to adaptively choose their goals (the grasping points, see Figure 5(c)) according to

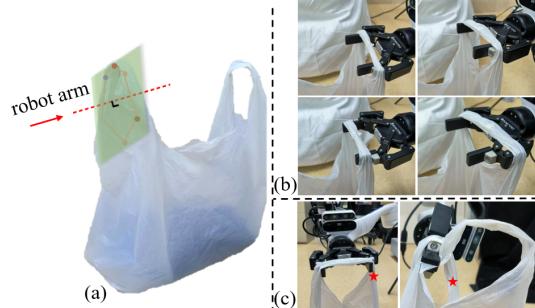


Figure 5: (a) Calculating the inserting direction.
(b) Different initial stages for tying Ian Knot.
(c) End-effector view before action 1 and action 2.
Red stars represent goals.

Table 1: Success rates of different methods. FSR of baseline methods are calculated by using our method to complete the missing parts of them.

Metrics	Plastic Bags	DI	RP	VBE-S	VBE-T	HKT	Ours
ASR	Red, Black, White	4.0%(2/60)	13.3%(8/60)	70.0%(42/60)	65.0%(39/60)	-	88.3%(53/60)
	Grey (Unlearned)	5.0%(1/20)	10.0%(2/20)	40.0%(8/20)	30.0%(6/20)	-	80.0%(16/20)
KSR	Red, Black, White	-	-	-	-	41.6%(25/60)	98.3%(59/60)
	Grey (Unlearned)	-	-	-	-	40.0%(8/20)	90.0%(18/20)
FSR	Red, Black, White	4.0%(2/60)	13.3%(8/60)	70.0%(42/60)	63.3%(38/60)	38.3%(23/60)	86.7%(52/60)
	Grey (Unlearned)	0.0%(0/20)	10.0%(2/20)	30.0%(6/20)	25.0%(5/20)	25.0%(5/20)	80.0%(16/20)

265 actual situations. To this end, we train a 4-layer convolutional neural network using Huber loss
 266 with hidden size=256 to regress grasping points for both actions directly in \mathbb{R}^3 with images input,
 267 as shown in Figure 5(c). Following up demonstrations in 3.3, we extend each demonstration with
 268 a knot tying part. For each demonstration, we record two images from right and left end-effector
 269 views respectively for action 1 and action 2, and record corresponding grasping locations in \mathbb{R}^3 .

270 4 Experimental Results

271 4.1 Metrics and Baselines

272 We evaluate our method on success rates of knot tying on three plastic bags that have demonstra-
 273 tions on them and a new type of plastic bag (grey) that has no demonstration on it to show the
 274 generalization ability of our method. Concretely, we evaluate: a) Adjusting Success Rate (**ASR**):
 275 if the plastic bag is successfully adjusted to a standing pose from a randomly dropped initial state
 276 to facilitate a valid inserting and lifting action (evaluated by human); b) Knot Tying Success Rate
 277 (**KSR**): if the plastic bag is successfully knotted from a randomly in-air initial stage. c) Full Task
 278 Success Rate (**FSR**): if the plastic bag is successfully knotted from a randomly dropped initial state.
 279 For ASR and FSR, we say one attempt fails if it does not succeed in 10 steps of grasping. We show
 280 the effectiveness of different modules in our method by a set of ablated versions:

281 **Directly Inserting without Adjusting (DI):** This method aims to directly find the handles of the
 282 plastic bag to lift it with the same method in 3.4 when it is just dropped from the air by the keypoint
 283 skeleton detected at the initial state. This is used to show the necessity of IIM.

284 **Random Picking (RP):** The robot randomly picks up a point of the plastic bag based on the seg-
 285 mented plastic bag area from top-down view image I_{right} and randomly generated a single grasping
 286 point in this area. This is used to show the effectiveness of the spatial action maps ϕ .

287 **Vision Based Keypoint Skeleton Evaluation (VBE):** This method evaluates the goodness of the
 288 adjusted pose of the plastic bag directly from images rather than using a GCN to process the keypoint
 289 skeleton. This is used to show the effectiveness of the task progress module G . We use channel-wise
 290 concatenated side view images (**VBE-S**) and top-down view images (**VBE-T**) for classifying stages
 291 in this baseline respectively.

292 **Hard-coded Knot Tying (HKT):** This method ties the Ian Knot with the same primitives of our
 293 method but use hard-coded goals for each action. This is used to show the effectiveness of learned
 294 motion primitives by CNN.

295 4.2 Quantitative Results

296 Table 1 shows the ASR, KSR, and FSR of different baselines. Our method achieves the best success
 297 rates in all metrics. For ASR, DI can barely knot the plastic bag, since most of the initial states do
 298 not support high-confidence keypoint detection (as shown in Figure 6), but occasionally a randomly
 299 dropped plastic bag can form a good standing pose. RP achieves similar results. Random picking
 300 can stretch out the plastic bag to some extent, but it can never end up with a good standing pose,
 301 which needs a dual-arm picking action p_2 . VBE methods achieve much better results than the above
 302 two methods on the training plastic bags, especially for VBE-S. However, their performance drops
 303 dramatically on the grey plastic bag. This is because although images can also provide information
 304 about poses of plastic bags, they do not extract and use the essential information for evaluating poses
 305 such as the keypoint skeleton, and are susceptible to different illumination conditions. Meanwhile
 306 we only provide hundreds of images for each plastic bag, which may not be sufficient for image-
 307 based classification. That is why VBE methods are not as good as IIM on training plastic bags and
 308 lack generalization abilities to new plastic bags. For KSR, the hard-coded knot tying method can
 309 achieve an average 41.25% success rate. Most of failures of HKT happen in action 2. The goal
 310 position in action 2 varies because it is affected by the grasping results of action 1. A hard coded
 311 action 2 will miss the front side of the right handle or just grasp a wrong part of the plastic bag.

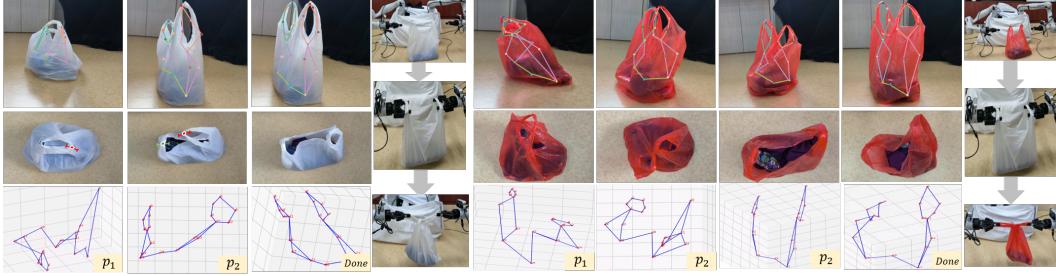


Figure 6: IIM processes in white and red plastic bags. Top row: side view and detected keypoint skeleton of plastic bags at each step. Middle row: top-down view and the predicted grasping points and directions. Bottom row: reconstructed 3D keypoint skeletons and the predicted actions at each step. We can see at both initial states, handles interleave with each other, but the robot can finally adjust them to standing poses with iterative pick-and-place actions.

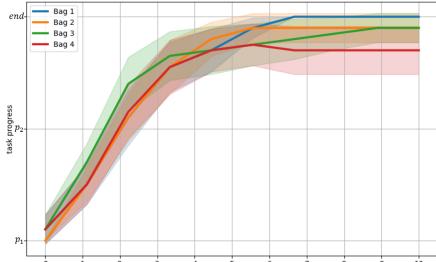


Figure 7: All task progresses along with grasping steps.

312 4.3 Qualitative Results

313 We show the qualitative results of white and red plastic bags in Figure 6 and all task progress curves
 314 evaluated by G (mean and std) in Figure 7. IIM shows great generalization ability to adjust plastic
 315 bags from different initial states to the standing poses. Most IIM processes finish in 4 steps. Some
 316 of them have more steps to finish because of harder initial configurations: one handle is too closer
 317 to the surface below it, or two handles interweave. In this situation the robot may grasp up more
 318 parts than one handle, but after several pick-and-place, the interlaced part can disentangle from each
 319 other and subsequent graspings become normal, as shown in the red plastic bag in Figure 6.

320 **Failure Cases:** The inaccuracy of depth estimation is the main reason that causes a failed grasping.
 321 Some graspings fail to catch the bag (above the right grasping point). Others exceed the right
 322 grasping depth. The first case will make the robot stuck in a loop. The second case will make the
 323 robot grasp the body of the plastic bag directly. This may lead to a violent jolt when placing the
 324 plastic bag, which almost equals starting all over again. Sometimes the filler items fall out, which
 325 leads to a failed grasping, as shown in Figure 8(c). Most of the attempts can finish in 5 grasps, but
 326 some attempts (around 15%) need more steps due to the above reason, and 13.3% of attempts fail.

327 5 Discussion and Future Works

328 In this work we propose Iterative Interactive Modeling for knotting plastic bags that are randomly
 329 dropped from the air with a dual arm robot. We show the effectiveness of various visual learning
 330 methods such as spatial-action maps and keypoint detection models on plastic bags. IIM is a
 331 general type of interactive perception [52] for modeling complex elasto-plastic objects: interacting with
 332 objects to complete an explicit representation model. Modules in our method can be replaced by
 333 other specific methods for different objects. For example, the representation model (keypoint skele-
 334 ton in this work) can be dense descriptors [8, 27, 28] or partical-based graphs [13, 14, 36], and the
 335 completion algorithm (imitation learning in this work) can be graph-based completion algorithms.

336 The limitations of our work are: 1) The keypoint detector depends on our custom dataset *PBpose*,
 337 and IIM is based on human demonstrations. These make our method lack quick extensibility to
 338 other objects and situations. 2) The assumption of this work is that the plastic bag can stand still
 339 after being adjusted. For some very soft/hard plastic bags, or some fillers that do not support stand-
 340 ing, our method would not be effective. 3) Our method cannot handle extremely difficult initial
 341 configurations such as the handles being pressed underneath the plastic bag. Future works may seek
 342 to employ more actions such as pushing to tackle these situations.

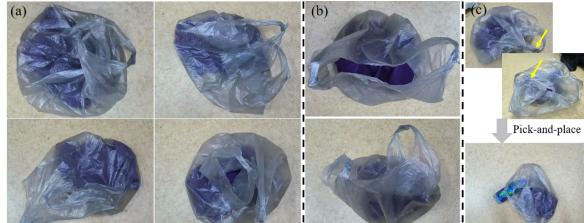


Figure 8: (a) Different initial configurations. (b) Special initial configurations that support directly lifting. (c) Failure cases. The filler is thrown out.

343 **References**

- 344 [1] A. Nair, D. Chen, P. Agrawal, P. Isola, P. Abbeel, J. Malik, and S. Levine. Combining self-
345 supervised learning and imitation for vision-based rope manipulation. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 2146–2153. IEEE, 2017.
- 346
- 347 [2] S. Kudoh, T. Gomi, R. Katano, T. Tomizawa, and T. Suehiro. In-air knotting of rope by a dual-
348 arm multi-finger robot. In *2015 IEEE/RSJ International Conference on Intelligent Robots and*
349 *Systems (IROS)*, pages 6202–6207. IEEE, 2015.
- 350 [3] H. Ha and S. Song. Flingbot: The unreasonable effectiveness of dynamic manipulation for
351 cloth unfolding. In *Conference on Robot Learning*, pages 24–33. PMLR, 2022.
- 352 [4] R. Lee, D. Ward, A. Cosgun, V. Dasagi, P. Corke, and J. Leitner. Learning arbitrary-goal fabric
353 folding with one hour of real robot experience. *arXiv preprint arXiv:2010.03209*, 2020.
- 354 [5] A. Namiki and S. Yokosawa. Robotic origami folding with dynamic motion primitives. In
355 *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages
356 5623–5628. IEEE, 2015.
- 357 [6] S. Duenser, J. M. Bern, R. Poranne, and S. Coros. Interactive robotic manipulation of elastic
358 objects. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*,
359 pages 3476–3481. IEEE, 2018.
- 360 [7] B. Shen, Z. Jiang, C. Choy, L. J. Guibas, S. Savarese, A. Anandkumar, and Y. Zhu. Acid:
361 Action-conditional implicit visual dynamics for deformable object manipulation. *arXiv preprint arXiv:2203.06856*, 2022.
- 362
- 363 [8] P. R. Florence, L. Manuelli, and R. Tedrake. Dense object nets: Learning dense visual object
364 descriptors by and for robotic manipulation. *arXiv preprint arXiv:1806.08756*, 2018.
- 365 [9] K. Suzuki, M. Kanamura, Y. Suga, H. Mori, and T. Ogata. In-air knotting of rope using dual-
366 arm robot based on deep learning. In *2021 IEEE/RSJ International Conference on Intelligent*
367 *Robots and Systems (IROS)*, pages 6724–6731. IEEE, 2021.
- 368 [10] A. Wang, T. Kurutach, K. Liu, P. Abbeel, and A. Tamar. Learning robotic manipulation through
369 visual planning and acting. *arXiv preprint arXiv:1905.04411*, 2019.
- 370 [11] Y. Wu, W. Yan, T. Kurutach, L. Pinto, and P. Abbeel. Learning to manipulate deformable
371 objects without demonstrations. *arXiv preprint arXiv:1910.13439*, 2019.
- 372 [12] X. Lin, Y. Wang, Z. Huang, and D. Held. Learning visible connectivity dynamics for cloth
373 smoothing. In *Conference on Robot Learning*, pages 256–266. PMLR, 2022.
- 374 [13] H. Shi, H. Xu, Z. Huang, Y. Li, and J. Wu. Robocraft: Learning to see, simulate, and shape
375 elasto-plastic objects with graph networks. *arXiv preprint arXiv:2205.02909*, 2022.
- 376 [14] Y. Li, J. Wu, R. Tedrake, J. B. Tenenbaum, and A. Torralba. Learning particle dynamics for
377 manipulating rigid bodies, deformable objects, and fluids. *arXiv preprint arXiv:1810.01566*,
378 2018.
- 379 [15] Z. Xu, Z. He, J. Wu, and S. Song. Learning 3d dynamic scene representations for robot
380 manipulation. *arXiv preprint arXiv:2011.01968*, 2020.
- 381 [16] L. Manuelli, W. Gao, P. Florence, and R. Tedrake. kpam: Keypoint affordances for category-
382 level robotic manipulation. *arXiv preprint arXiv:1903.06684*, 2019.
- 383 [17] A. Zeng, S. Song, J. Lee, A. Rodriguez, and T. Funkhouser. Tossingbot: Learning to throw
384 arbitrary objects with residual physics. *IEEE Transactions on Robotics*, 36(4):1307–1319,
385 2020.
- 386 [18] A. Zeng, S. Song, S. Welker, J. Lee, A. Rodriguez, and T. Funkhouser. Learning synergies
387 between pushing and grasping with pervised deep reinforcement learning. In *2018 IEEE/RSJ*
388 *International Conference on Intelligent Robots and Systems (IROS)*, pages 4238–4245. IEEE,
389 2018.

- 390 [19] J. Smolen and A. Patriciu. Deformation planning for robotic soft tissue manipulation. In
 391 *2009 Second International Conferences on Advances in Computer-Human Interactions*, pages
 392 199–204. IEEE, 2009.
- 393 [20] S. Luo, J. Bimbo, R. Dahiya, and H. Liu. Robotic tactile perception of object properties: A
 394 review. *Mechatronics*, 48:54–67, 2017.
- 395 [21] J. Wu, X. Sun, A. Zeng, S. Song, J. Lee, S. Rusinkiewicz, and T. Funkhouser. Spatial action
 396 maps for mobile manipulation. *arXiv preprint arXiv:2004.09141*, 2020.
- 397 [22] J. Li, S. Bian, A. Zeng, C. Wang, B. Pang, W. Liu, and C. Lu. Human pose regression with
 398 residual log-likelihood estimation. In *Proceedings of the IEEE/CVF International Conference*
 399 *on Computer Vision*, pages 11025–11034, 2021.
- 400 [23] J. Maitin-Shepard, M. Cusumano-Towner, J. Lei, and P. Abbeel. Cloth grasp point detection
 401 based on multiple-view geometric cues with application to robotic towel folding. In *2010 IEEE*
 402 *International Conference on Robotics and Automation*, pages 2308–2315. IEEE, 2010.
- 403 [24] D. Seita, N. Jamali, M. Laskey, R. Berenstein, A. K. Tanwani, P. Baskaran, S. Iba, J. Canny, and
 404 K. Goldberg. Robot bed-making: Deep transfer learning using depth sensing of deformable
 405 fabric. *arXiv preprint arXiv:1809.09810*, 26, 2018.
- 406 [25] S. Miller, J. Van Den Berg, M. Fritz, T. Darrell, K. Goldberg, and P. Abbeel. A geometric
 407 approach to robotic laundry folding. *The International Journal of Robotics Research*, 31(2):
 408 249–267, 2012.
- 409 [26] T. D. Kulkarni, A. Gupta, C. Ionescu, S. Borgeaud, M. Reynolds, A. Zisserman, and V. Mnih.
 410 Unsupervised learning of object keypoints for perception and control. *Advances in neural*
 411 *information processing systems*, 32, 2019.
- 412 [27] A. Ganapathi, P. Sundaresan, B. Thananjeyan, A. Balakrishna, D. Seita, J. Grannen, M. Hwang,
 413 R. Hoque, J. E. Gonzalez, N. Jamali, et al. Learning dense visual correspondences in simulation
 414 to smooth and fold real fabrics. In *2021 IEEE International Conference on Robotics and*
 415 *Automation (ICRA)*, pages 11515–11522. IEEE, 2021.
- 416 [28] P. Florence, L. Manuelli, and R. Tedrake. Self-supervised correspondence in visuomotor policy
 417 learning. *IEEE Robotics and Automation Letters*, 5(2):492–499, 2019.
- 418 [29] L. Yen-Chen, P. Florence, J. T. Barron, T.-Y. Lin, A. Rodriguez, and P. Isola. Nerf-
 419 supervision: Learning dense object descriptors from neural radiance fields. *arXiv preprint*
 420 *arXiv:2203.01913*, 2022.
- 421 [30] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese. 3d-r2n2: A unified approach for single
 422 and multi-view 3d object reconstruction. In *European conference on computer vision*, pages
 423 628–644. Springer, 2016.
- 424 [31] T. Groueix, M. Fisher, V. G. Kim, B. C. Russell, and M. Aubry. A papier-mâché approach to
 425 learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and*
 426 *pattern recognition*, pages 216–224, 2018.
- 427 [32] B. Deng, K. Genova, S. Yazdani, S. Bouaziz, G. Hinton, and A. Tagliasacchi. Cvxnet: Learn-
 428 able convex decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
 429 *and Pattern Recognition*, pages 31–44, 2020.
- 430 [33] C. Jiang, J. Huang, A. Tagliasacchi, L. Guibas, et al. Shapeflow: Learnable deformations
 431 among 3d shapes. *arXiv preprint arXiv:2006.07982*, 2020.
- 432 [34] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla.
 433 Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International*
 434 *Conference on Computer Vision*, pages 5865–5874, 2021.
- 435 [35] D. Mrowca, C. Zhuang, E. Wang, N. Haber, L. F. Fei-Fei, J. Tenenbaum, and D. L. Yamins.
 436 Flexible neural representation for physics prediction. *Advances in neural information process-
 437 ing systems*, 31, 2018.

- 438 [36] Y. Li, T. Lin, K. Yi, D. Bear, D. Yamins, J. Wu, J. Tenenbaum, and A. Torralba. Visual
 439 grounding of learned physical models. In *International conference on machine learning*, pages
 440 5927–5936. PMLR, 2020.
- 441 [37] Y. Hu, L. Anderson, T.-M. Li, Q. Sun, N. Carr, J. Ragan-Kelley, and F. Durand. Diffitaichi:
 442 Differentiable programming for physical simulation. *arXiv preprint arXiv:1910.00935*, 2019.
- 443 [38] P. Holl, V. Koltun, and N. Thuerey. Learning to control pdes with differentiable physics. *arXiv
 444 preprint arXiv:2001.07457*, 2020.
- 445 [39] Y. Li, S. Li, V. Sitzmann, P. Agrawal, and A. Torralba. 3d neural scene representations for
 446 visuomotor control. In *Conference on Robot Learning*, pages 112–123. PMLR, 2022.
- 447 [40] B. Frank, C. Stachniss, N. Abdo, and W. Burgard. Efficient motion planning for manipulation
 448 robots in environments with deformable objects. In *2011 IEEE/RSJ International Conference
 449 on Intelligent Robots and Systems*, pages 2180–2185. IEEE, 2011.
- 450 [41] M. Saha and P. Isto. Manipulation planning for deformable linear objects. *IEEE Transactions
 451 on Robotics*, 23(6):1141–1150, 2007.
- 452 [42] M. Moll and L. E. Kavraki. Path planning for deformable linear objects. *IEEE Transactions
 453 on Robotics*, 22(4):625–636, 2006.
- 454 [43] D. McConachie, M. Ruan, and D. Berenson. Interleaving planning and control for deformable
 455 object manipulation. In *Robotics Research*, pages 1019–1036. Springer, 2020.
- 456 [44] D. Berenson. Manipulation of deformable objects without modeling and simulating deforma-
 457 tion. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages
 458 4525–4532. IEEE, 2013.
- 459 [45] Z. Hu, P. Sun, and J. Pan. Three-dimensional deformable object manipulation using fast online
 460 gaussian process regression. *IEEE Robotics and Automation Letters*, 3(2):979–986, 2018.
- 461 [46] C. Wang, S. Wang, B. Romero, F. Veiga, and E. Adelson. Swingbot: Learning physical fea-
 462 tures from in-hand tactile exploration for dynamic swing-up manipulation. In *2020 IEEE/RSJ
 463 International Conference on Intelligent Robots and Systems (IROS)*, pages 5633–5640. IEEE,
 464 2020.
- 465 [47] V. Chu, I. McMahon, L. Riano, C. G. McDonald, Q. He, J. M. Perez-Tejada, M. Arrigo,
 466 T. Darrell, and K. J. Kuchenbecker. Robotic learning of haptic adjectives through physical
 467 interaction. *Robotics and Autonomous Systems*, 63:279–292, 2015.
- 468 [48] A. X. Lee, A. Gupta, H. Lu, S. Levine, and P. Abbeel. Learning from multiple demonstra-
 469 tions using trajectory-aware non-rigid registration with applications to deformable object manipula-
 470 tion. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*,
 471 pages 5265–5272. IEEE, 2015.
- 472 [49] C. Matl and R. Bajcsy. Deformable elasto-plastic object shaping using an elastic hand and
 473 model-based reinforcement learning. In *2021 IEEE/RSJ International Conference on Intelli-
 474 gent Robots and Systems (IROS)*, pages 3955–3962. IEEE, 2021.
- 475 [50] J. Matas, S. James, and A. J. Davison. Sim-to-real reinforcement learning for deformable
 476 object manipulation. In *Conference on Robot Learning*, pages 734–743. PMLR, 2018.
- 477 [51] J. L. Schonberger and J.-M. Frahm. Structure-from-motion revisited. In *Proceedings of the
 478 IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.
- 479 [52] J. Bohg, K. Hausman, B. Sankaran, O. Brock, D. Kragic, S. Schaal, and G. S. Sukhatme. Inter-
 480 active perception: Leveraging action in perception and perception in action. *IEEE Transactions
 481 on Robotics*, 33(6):1273–1291, 2017.