

Vegetable Peeling: A Case Study in Constrained Dexterous Manipulation

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Recent studies have made significant progress in addressing dexterous
2 manipulation problems, particularly in in-hand object reorientation. However,
3 there are few existing works that explore the potential utilization of developed
4 dexterous manipulation controllers for downstream tasks. In this study, we focus
5 on constrained dexterous manipulation for food peeling. Food peeling presents
6 various constraints on the reorientation controller, such as the requirement for the
7 hand to securely hold the object after reorientation for peeling. We propose a
8 simple system for learning a reorientation controller that facilitates the subsequent
9 peeling task. Videos are available at: [https://sites.google.com/view/
10 food-peeling](https://sites.google.com/view/food-peeling).

11 **Keywords:** Dexterous manipulation, In-hand object reorientation, vegetable peel-
12 ing

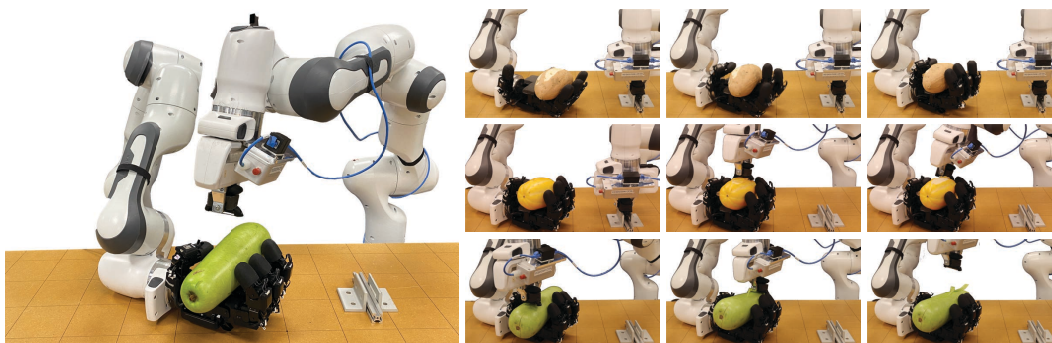


Figure 1: We present a dexterous manipulation system that utilizes an Allegro hand mounted on a Franka robot arm to reorient food items for downstream peeling. The other Franka robot arm uses its gripper to grasp a peeler for peeling. The reorientation controller for the Allegro hand is learned through reinforcement learning, while the peeling is performed via teleoperation. On the left of the figure, we show the whole system. On the right, from the top row to the bottom row, we illustrate the processes of reorienting a sweet potato, and peeling a melon and a squash.

13 1 INTRODUCTION

14 Having robots perform food preparation tasks has been of great interest in robotics. Imagine the
15 scenario of making mashed potatoes, where a critical step is to peel potatoes. Humans peel potatoes
16 by grasping the potato in one hand and using the second hand to actuate a peeler to remove the
17 potato’s skin. After a part of the potato is peeled, it is rotated while being held in the hand (i.e.,
18 *in-hand manipulation*) and peeled again. The sequence of rotating and peeling continues until all of
19 the potato’s skin is removed. In this work, we present a robotic system that can re-orient different
20 vegetables using an Allegro hand in a way that their skin can be peeled using another manipulator.
21 Our setup is shown in Figure 1.

22 In-hand rotation of vegetables is an instance of dexterous manipulation problem [1], a family of
23 tasks that involves continuously controlling the force on an object while it is moving with respect
24 to the fingertips [2, 3]. The challenges in dexterous manipulation stem from the frequent making
25 and breaking of contact, issues in contact modeling, high-dimensional control space, perception
26 challenges due to severe occlusions, etc. A body of work made simplifying assumptions such as
27 manipulating convex objects [4, 5, 1, 6], small finger motions [7, 8, 9], slow or quasi-static motion or
28 manipulating a few specific objects [10, 7, 8] to leverage trajectory optimization or planning-based
29 methods to achieve in-hand object re-orientation [1, 7, 8, 9, 6, 4, 5, 10]. Another line of work has
30 used reinforcement learning for in-hand re-orientation [11, 12, 13, 14, 15] and recent works have
31 leveraged simulation to train policies capable of dynamically re-orienting a diverse set of new objects
32 in real-time and in the real world [11, 12].

33 There are several challenges in adapting re-orientation controllers for a downstream task such
34 as peeling vegetables. These challenges stem from the fact that controllers optimized for re-
35 orientation [16, 13, 14, 15, 12] are only optimized to continuously reorient the object and not
36 to satisfy numerous constraints arising from task-specific requirements. For instance, peeling vegeta-
37 bles requires the hand to **first stop** re-orienting the object and then for the peeler to peel the vegetable.
38 Many prior works solve a version of the re-orientation problem where the object is continuously
39 rotated [17, 16, 13] or otherwise perform quasistatic re-orientation [8]. Stopping and re-starting
40 *dynamic* re-orientation is difficult due to the challenge of dealing with the object’s inertia. **Second**,
41 the hand needs to *hold* the object firmly enough to resist forces applied by the peeler. The closest
42 work that attempts to hold the object at a target configuration [12] is only able to loosely hold the
43 object which is insufficient for resisting forces. **Third**, the hand needs to reorient the vegetable *along*
44 *a specific axis in place*. Here, the specific axis refers to the rotational axis on the object that is parallel
45 to the peeling direction. Similar to how humans reorient vegetables for peeling, it is desirable for the
46 hand to reorient the object in place so that multiple consecutive cycles of reorientation and peeling
47 can be performed. If the object substantially shifts its position during reorientation, the controller
48 will struggle to reorient and hold the object at future time steps. **Fourth**, when the vegetable is held
49 stationary the fingers should *not obstruct* the top surface of the vegetable to ensure that the peeler can
50 peel the vegetable.

51 While in-hand object reorientation has been widely studied [11, 12, 16, 18, 13, 17], no prior works
52 can satisfy the constraints mentioned above. Yet, these constraints become critical for downstream
53 dexterous manipulation beyond object re-orientation. We use vegetable peeling as a case study to
54 investigate the challenges and solutions for building a dexterous manipulation system that can operate
55 under constraints. We develop a framework where we leverage reinforcement learning in simulation
56 to train a policy that can perform object re-orientation under constraints. For the peeling, we have a
57 human teleoperate the peeler. Our contributions are as follows:

- 58 1. A framework for solving dexterous manipulation problems under constraints.
- 59 2. We propose a method that can make RL policy learn to stop its motion and hold objects
60 firmly in hand – a critical behavior for many downstream dexterous manipulation problems.
- 61 3. We present a step towards a robotic system capable of peeling diverse vegetables with differ-
62 ent shapes, masses, and material properties while holding and manipulating the vegetables
63 in hand.

64 2 RELATED WORK

65 **In-hand Object Reorientation:** Dexterous manipulation involves the use of high degrees-of-freedom
66 (DoF) manipulators for object manipulation [19]. Its requirement for high-dimensional real-time
67 control and its nature of frequent contact-making and breaking present grand challenges to roboticists.
68 Recently, there has been a growth of interest in a particular instance of dexterous manipulation
69 problems: in-hand object reorientation. This problem is of particular interest as it is a necessary
70 step in many tool-use scenarios. For example, to use a screwdriver for tightening a screw, one has
71 to reorient the screwdriver to align it with the screw. We can cluster the works in in-hand object

72 reorientation from many aspects. For example, from the perspective of sensory information, [20]
 73 studies open-loop cube reorientation without using any sensors, [21, 5, 16, 10, 22] use motion capture
 74 system or special tracking markers for object reorientation, [17] uses proprioceptive sensors such as
 75 joint encoders, [23, 24, 15, 14] use tactile sensors and [25, 16, 12, 18] utilize vision sensors. In terms
 76 of the dynamics of the system, [7, 8, 9] achieved object reorientation under the assumption of quasi-
 77 static motion where object moves slowly and its inertia effect can be ignored, while [15, 16, 12, 14, 26]
 78 focuses on dynamic object reorientation where object is manipulated in a fast and dynamic way. To
 79 make in-hand object manipulation useful for downstream tool use tasks, one important aspect of the
 80 skill is the ability of stably and firmly holding the object in end of the policy rollout. While many
 81 prior works on dynamic manipulation such as [16, 10, 14, 15, 17] only consider endlessly rotating
 82 the object in hand and cannot stop the object stably when the object reaches the goal orientation,
 83 some works such as [12, 26] try to develop controllers that can reorient objects in hand and also hold
 84 the object in the goal orientation. Our work studies dynamic in-hand object manipulation with the
 85 capability of stopping objects stably in hand.

86 **Reinforcement Learning for Contact-rich Tasks:** Contact-rich tasks are particularly challenging
 87 due to the difficulty in modeling the system dynamics, especially when the tasks are performed in the
 88 wild, outside of a constrained and controlled setting. Examples of such tasks include quadruped robots
 89 hiking in mountains and robot hands reorienting various everyday objects. There have been many
 90 works using reinforcement learning to learn controllers for solving contact-rich tasks [27, 16, 13, 28,
 91 29, 30]. In the real world, robots typically only have access to a limited amount of state information of
 92 the system due to the lack of sensors or the challenges in setting up the sensors. Using reinforcement
 93 learning to learn controllers from scratch with limited sensory information tends to be data-inefficient.
 94 One way to speed up policy learning is to provide asymmetric information to the policy and value
 95 function, where the value function observes much more privileged information [16, 13, 27, 31].
 96 Another method is to decouple policy learning into two stages: a reinforcement learning stage where
 97 agents (teacher) observe privileged fully-observable state information, and an imitation learning stage
 98 where the policy with limited sensory observation input (student) learns to imitate the policy with
 99 fully-observable state information. This approach has been successfully applied to various contact-
 100 rich problems such as locomotion [32, 33, 30, 34, 35] and dexterous manipulation [11, 12, 17].
 101 Our pipeline is built upon the idea of teacher-student policy learning and has made several key
 102 improvements, which we will detail below.

103 3 METHOD

104 Peeling requires a reorientation controller that can stop its motion and firmly hold objects after
 105 reorientation. The first step in stopping is to decide when re-orientation should be stopped. One
 106 possibility is to have a perception system predict the desired rotation angle after which the next round
 107 of peeling would be performed. To accomplish the goal, the robot would need to track changes in
 108 object pose and compare it with the target rotation angle. However, accurately estimating object pose
 109 is challenging, especially when generalization to new objects is necessary [36, 16, 13].

110 One of our insights is that instead of training a predictor for desired rotation angle and object pose
 111 estimation, it can be *easier* and *sufficient* to train a *binary vision classifier* that detects in real-time
 112 when the peeled part has been turned over. With such a classifier, the reorientation controller’s job is
 113 simply to keep reorienting the object until it receives a stop signal. In this formulation, unlike prior
 114 works [11, 12], the reorientation controller is not conditioned on target orientation but rather on a
 115 stop signal. Formally, the policy takes as input a binary variable $I_t^{stop} \in \{0, 1\}$ representing the stop
 116 signal. If $I_t^{stop} = 1$, the policy should stop immediately and ensure the fingers stably and firmly hold
 117 the object. Otherwise, the policy should continue reorienting the object. Note that in this work, we
 118 focus on learning the reorientation controller, leaving integration of a vision classifier to future work.

119 The next question is how to train such a policy. Using RL to train the policy from scratch can be
 120 challenging and requires extensive reward shaping because $I_t^{stop} = 1$ is a rare event in an episode,

121 and when the I_t^{stop} is flipped to one from zero, the policy needs to quickly stop the motion posing a
 122 hard-exploration challenge.

123 Prior works [11, 12] show success in training a goal-conditioned object reorientation controller. Can
 124 we leverage a goal-conditioned reorientation controller to train a controller that reacts to a stop signal?
 125 It turns out we can formulate this using the teacher-student learning framework [11, 12, 37, 34, 33].
 126 Specifically, we can use RL to train a goal-conditioned controller that reorients an object by random
 127 goal angles along its rotational axis. This acts as the teacher. Next, we can use imitation learning
 128 (specifically DAGGER [38]) to train a controller conditioned on the stop signal to imitate the teacher.
 129 The stop signal can be generated during training by checking if the orientation distance to the goal is
 130 below a threshold. Using imitation learning bypasses the hard exploration challenge.

131 3.1 Teacher Policy Learning: Reorient and Stop

132 We train the teacher policy to re-orient the object along a pre-defined axis and stop (see Figure 2a).
 133 The teacher is formulated as a goal-conditioned policy $\mathbf{a}_t^{\mathcal{E}} = \pi^{\mathcal{E}}(\mathbf{o}_t^{\mathcal{E}}, \mathbf{a}_{t-1}, g)$, where \mathcal{E} represents
 134 variables for the teacher policy, \mathbf{o}_t is the observation, \mathbf{a}_t is the action command, g is the goal
 135 representing the amount by which the object needs to be re-oriented. g is randomly and uniformly
 136 sampled from $[1.57, 4.0]$ rad during training.

137 While the teacher policy’s formulation is similar to that in prior works [11, 12], we propose (i) a
 138 much simpler reward function, (ii) new success criteria that effectively encourages the policy to stop
 139 the object and firmly hold it, and (iii) an interpolation scheme that enables smoother policy actions in
 140 the real world.

141 3.1.1 Observation and Action Space

142 $\mathbf{o}_t^{\mathcal{E}}$ includes joint positions and velocities, the fingertip poses and velocities, object pose and velocity,
 143 the distance between the current object orientation and the goal orientation, and whether any of the
 144 fingertips touch the object. \mathbf{a}_t is the delta joint position command. The neural network policy runs at
 145 12Hz.

146 3.1.2 Reward Function

147 For the task of in-hand re-orientation, we found a simple way to specify the reward function.
 148 Specifically, we manually move the real Allegro hand to a good pose where the constraints mentioned
 149 above are satisfied (e.g., the fingers do not cover the food item), and the fingers touch the object and
 150 are ready to reorient it. We record the joint positions as \mathbf{q}^{demo} . During training in simulation, we
 151 encourage the joint positions at any time step to be close to \mathbf{q}^{demo} .

152 Overall, our reward function is as follows:

$$r_t = c_1 \mathbb{1}(\text{Task successful}) + c_2 \frac{1}{|\Delta\theta_t| + \epsilon_\theta} \quad (1)$$

$$+ c_3 \|\mathbf{q}_t - \mathbf{q}^{demo}\|_2^2 \quad (2)$$

153 where c_1, c_2, c_3 are coefficients. $\Delta\theta_t$ is the distance between the object’s current and goal orientation.
 154 The first two terms are task rewards for object reorientation. The last term is to regulate hand behavior.

155 3.1.3 Success Criteria

156 In a goal-conditioned object reorientation, a common way to claim the task successful is by checking
 157 if the distance between the object’s current and the goal orientation is smaller than a threshold value
 158 (orientation criterion $C_{ori} = \Delta\theta < \bar{\theta}$) [16, 13]. Another criterion is that all the fingertips should
 159 make contact with the object (contact criterion $C_{contact}$), a pre-requisite for firmly holding the object
 160 after reorientation. However, only checking these two criteria is insufficient to ensure the policy
 161 learns to stop the motion and hold the object firmly around the goal orientation, as discussed in [12].
 162 The policy can oscillate around the goal state due to observation and control delay and noise.

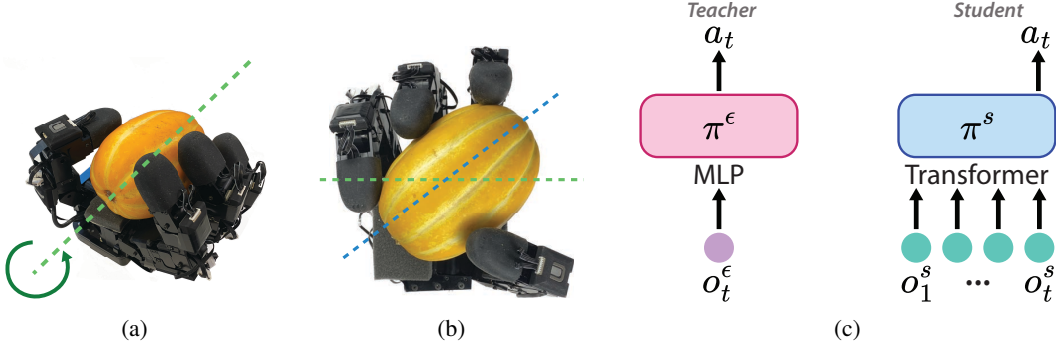


Figure 2: **(a)** shows an example of the rotational axis of a melon. **(b)** shows an example where the object’s orientation (the blue line) has a large deviation from the desired rotational axis (the green line). We reset the episode when this occurs. **(c)** shows the policy Architecture for the teacher and the student. In this figure, we use o_t to represent all the policy input at each time step.

163 To further encourage the policy to stop robot motion when the goal is reached and firmly hold the
 164 object, we propose adding time constraints to the success criteria: both C_{ori} and $C_{contact}$ should be
 165 continuously satisfied for \bar{T}^{succ} time steps. Adding this criterion makes the MDP partially observable
 166 since the policy’s observation lacks the knowledge of time. Therefore, to facilitate policy learning,
 167 we augment the observation space with a scalar indicator variable $I^{succ} = t^{succ} / \bar{T}^{succ} \in [0, 1]$,
 168 where t^{succ} is the number of consecutive steps satisfying C_{ori} and $C_{contact}$. The observation space
 169 becomes $o^\mathcal{E} := o^\mathcal{E} \oplus I^{succ}$. In this work, $\bar{\theta} = 0.2rad$, $\bar{T}^{succ} = 8$.

170 3.2 Student Policy Learning: Imitate and Stop

171 After learning a goal-conditional teacher policy $a_t^\mathcal{E} = \pi^\mathcal{E}(o_t^\mathcal{E}, a_{t-1}, g)$, the next question is how
 172 to train a real-world deployable student policy that can rotate the object in hand and hold it stably
 173 after reorientation. We propose conditioning the student policy on a stop signal $I_t^{stop} \in \{0, 1\}$:
 174 $a_t^S = \pi^S(o_t^S, a_{t-1}, I_t^{stop})$. In other words, the student policy should continue reorienting the object
 175 when $I_t^{stop} = 0$, but stably hold the object when $I_t^{stop} = 1$. This design choice provides flexibility in
 176 how we control the policy to stop the reorientation. For example, the policy could rotate the object
 177 for a pre-specified amount of time (i.e., set $I_t^{stop} = 1$ after t seconds). Alternatively, an external
 178 perception module could detect when the peeled part has fully turned over, triggering $I_t^{stop} = 1$ and
 179 the policy to stop the motion and hold the object immediately.

180 How can we use the learned goal-conditioned teacher policy to train a student policy that is conditioned
 181 on the stop signal? We can set the value for I_t^{stop} automatically during policy rollout based on the
 182 orientation distance $\Delta\theta_t$.

$$I_t^{stop} = \begin{cases} 0 & \text{if } \Delta\theta_t > \bar{\theta} \\ 1 & \text{otherwise} \end{cases}$$

183 3.2.1 Observation Space

184 In this work, we only use proprioceptive sensory information (joint positions q_t and velocities \dot{q}_t) as
 185 the observation input (o_t^S).

186 3.2.2 Policy Architecture

187 As the student policy only has access to a limited amount of sensory information (a POMDP setting),
 188 it is important to incorporate history information, as has been done in previous works [16, 13, 12].
 189 While [16, 13, 12] utilized RNNs to process history information, Transformers [39] have gained
 190 significant attention due to their improved performance and faster training in domains such as
 191 natural language processing. Therefore, in this work, we employ a Transformer-based policy
 192 architecture. $a_t^S = \pi^S(o_1^S, a_0, I_1^{stop}, \dots, o_t^S, a_{t-1}, I_t^{stop})$. The policy is a decoder-only attention

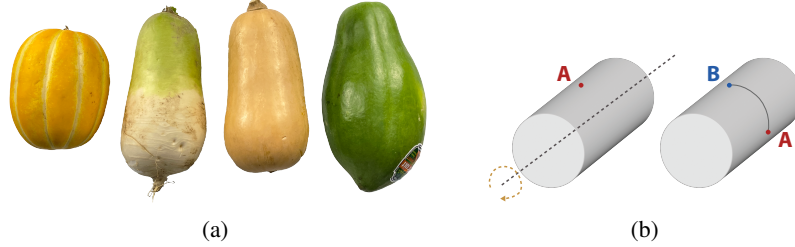


Figure 3: **(a)** shows the objects for evaluation: melon, radish, pumpkin, papaya. **(b)** shows the traveling distance. Before reorientation begins, we ensure a reference point (point A) is facing upward. After reorientation, we identify the point (point B) now facing upward. We then measure the distance from point A to point B along the contour.

193 network (Figure 2c) with three self-attention layers. The hidden size is 256, the intermediate size is
 194 512, and the number of attention heads is 8.

195 3.2.3 Training

196 The policy is trained using DAGGER [38].

197 3.3 Peeling via Teleoperation

198 We demonstrate that our reorientation controller can be used for downstream peeling tasks by
 199 teleoperating a Franka Panda robot arm to do the peeling. A 200 Hz operational space impedance
 200 controller [40] runs on the Panda arm, controlling for pose via torque, and an operator interacts with a
 201 Haption Virtuose™ 6D HF TAO¹ device. Bilateral position-position haptic coupling is done between
 202 the two devices. The controllers and haptic coupling are implemented using Drake [41].

203 4 RESULTS

204 To quantitatively evaluate the real-world policy transfer performance, we tested the controller on four
 205 vegetables (Figure 3a): a pumpkin (mass: 827g), a melon(623g), a radish(727g), a papaya(848g).

206 4.1 Traveling distance for a fixed amount of commanded motion time

207 The first question we want to answer is whether the learned policy can successfully reorient vegetables
 208 in the real world. In peeling, the width of the peeled part depends on the peeler’s width. Thus, it is
 209 more informative to measure how much the reorientation controller rotates an object by the traveling
 210 distance of a surface point, rather than the absolute rotation angle. Specifically, we mark a reference
 211 point P^{ref} on the object surface near the mid-point of its rotational axis. At the start, we ensure P^{ref}
 212 is centered and facing upward when held. After reorientation, we record the new point P^{new} that is
 213 now centered and facing upward. We then measure the contour length from P^{new} to P^{ref} along the
 214 surface (Figure 3b).

215 To demonstrate the capability of our controller to reorient real objects, we conducted two rounds of
 216 testing. Our controller is trained to stop motion when it receives a stop signal. In the first round, we
 217 sent the stop signal 3.5 seconds after the controller started rotating. In the second round, we sent the
 218 stop signal 7 seconds after start. We repeated each test 10 times. As shown in Figure 4a, the controller
 219 successfully reoriented all four food items by a sufficient amount for peeling. When commanded to
 220 reorient for 3.5s, 90% of tests reoriented the objects by at least 4cm. With 7s, 90% of tests reoriented
 221 objects by at least 7.3cm. Given more time, the controller reoriented objects by a larger amount.

¹<https://www.haption.com/en/products-en/virtuose-6d-tao-en.html#fa-download-downloads>

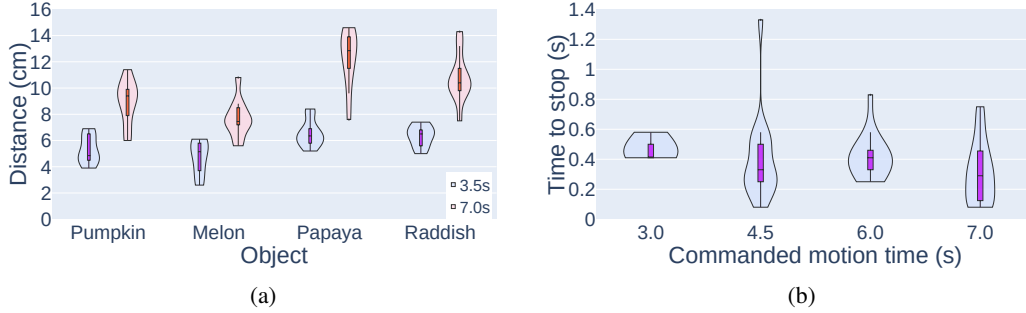


Figure 4: **(a)**: Violin plots showing the distribution of the traveling distance of a point on the object surface after the controller is commanded to rotate the object for 3.5 s and 7 s, respectively. **(b)**: Violin plot showing the distribution of time taken by the controller to transition from rotating the object in hand to firmly holding the object after receiving the stop signal. The x -axis represents the timing of the stop signal sent to the controller after it starts.

Table 1: Successful lifting rate (10 tests each)

Commanded motion time	Pumpkin	Melon	Papaya	Radish
3.5s	80%	90%	80%	90%
7s	100%	90%	100%	90%

222 4.2 How well does the controller track the commanded motion time?

223 As discussed in Section 3, if our controller can quickly respond to a stop signal at any time step, it
 224 can be combined with a perception system that tracks peeling progress. Hence, we measured how
 225 long it takes to stop the hand and object motion after receiving the stop signal. As shown in Figure 4b,
 226 the motion stops after 0.4s on average after the controller receives the stop signal.

227 4.3 Firm grasp after reorientation

228 To enable downstream peeling, the reorientation controller must learn to firmly grasp the object after
 229 stopping finger motion. We tested this by checking if the Allegro hand and object could be lifted in
 230 the air for 3s by only lifting the object. Table 1 shows that across objects and commanded times, the
 231 controller firmly grasped objects in 90% of tests.

232 4.4 Real-world Peeling

233 We evaluated whether the reorientation controller could reorient food items to facilitate peeling
 234 (Figure 1). Testing showed that peeling applied substantial pulling forces on objects. However, in
 235 most cases, the hand maintained a firm enough grasp to enable successful peeling.

236 4.5 Ablation study

237 4.5.1 Demo term in Reward function

238 We proposed using a keyframe demonstration to ease reward shaping. To evaluate its effectiveness,
 239 we compared learning curves of the teacher policies trained with and without the $c_3 \|\mathbf{q}_t - \mathbf{q}^{demo}\|_2^2$
 240 reward term. As shown in Figure 5a, adding the keyframe substantially improved learning. Addition-
 241 ally, it demonstrates that mimicking the keyframe pose via a single reward term effectively reduces
 242 the reward-shaping burden.

243 4.5.2 Transformer vs RNN

244 Different from prior works [16, 13, 11, 12], our student policy uses a Transformer architecture instead
 245 of an RNN architecture. We compared the learning performance of a Transformer-based policy and

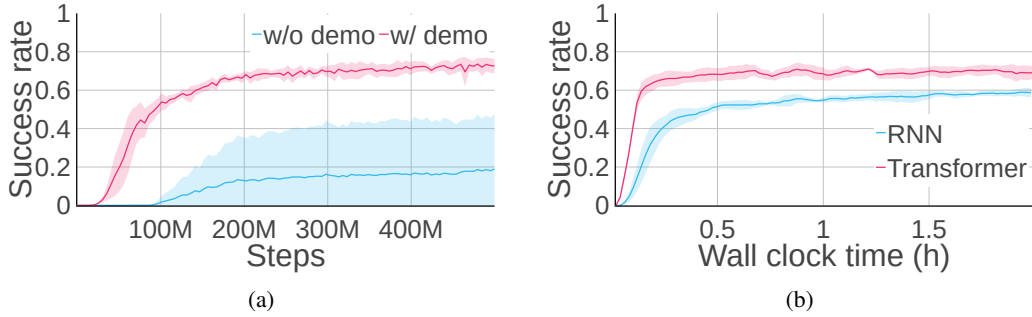


Figure 5: **(a)** shows learning curves of the teacher policies with or without $c_3 \|q_t - q^{demo}\|_2^2$ in the reward function. **(b)** shows the learning curve of student policies with a Transformer or RNN architecture with respect to the number of samples.

246 an RNN-based policy. Figure 5b shows that a Transformer-based policy learns much faster and gets
 247 better performance at convergence than an RNN-based policy.

248 5 DISCUSSIONS

249 The reorientation controller described in this study is a blind controller that relies solely on proprio-
 250 ceptive sensory information. Although it has shown the ability to successfully reorient heavy objects
 251 and securely hold them in place, performance could potentially be enhanced by incorporating visual
 252 and tactile feedback. For instance, visual information could help prevent objects from falling. Addi-
 253 tionally, future work could involve learning a peeling policy via behavior cloning on data collected
 254 via teleoperation to achieve full autonomy of the system.

References

- 255
- 256 [1] D. Rus. In-hand dexterous manipulation of piecewise-smooth 3-d objects. *The International*
257 *Journal of Robotics Research*, 18(4):355–381, 1999.
- 258 [2] M. T. Mason, J. K. Salisbury, and J. K. Parker. *Robot hands and the mechanics of manipulation*.
259 The MIT Press, 1989.
- 260 [3] N. C. Daffe, A. Rodriguez, R. Paolini, B. Tang, S. S. Srinivasa, M. Erdmann, M. T. Mason,
261 I. Lundberg, H. Staab, and T. Fuhlbrigge. Extrinsic dexterity: In-hand manipulation with
262 external forces. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*,
263 pages 1578–1585. IEEE, 2014.
- 264 [4] Y. Bai and C. K. Liu. Dexterous manipulation using both palm and fingers. In *2014 IEEE*
265 *International Conference on Robotics and Automation (ICRA)*, pages 1560–1565. IEEE, 2014.
- 266 [5] B. Sundaralingam and T. Hermans. Relaxed-rigidity constraints: kinematic trajectory optimiza-
267 tion and collision avoidance for in-grasp manipulation. *Autonomous Robots*, 43(2):469–483,
268 2019.
- 269 [6] I. Mordatch, Z. Popović, and E. Todorov. Contact-invariant optimization for hand manipulation.
270 In *Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer animation*,
271 pages 137–144, 2012.
- 272 [7] T. Pang, H. Suh, L. Yang, and R. Tedrake. Global planning for contact-rich manipulation via
273 local smoothing of quasi-dynamic contact models. *arXiv preprint arXiv:2206.10787*, 2022.
- 274 [8] A. S. Morgan, K. Hang, B. Wen, K. Bekris, and A. M. Dollar. Complex in-hand manipulation via
275 compliance-enabled finger gaiting and multi-modal planning. *IEEE Robotics and Automation*
276 *Letters*, 7(2):4821–4828, 2022.
- 277 [9] S. Abondance, C. B. Teeple, and R. J. Wood. A dexterous soft robotic hand for delicate in-hand
278 manipulation. *IEEE Robotics and Automation Letters*, 5(4):5502–5509, 2020.
- 279 [10] A. Nagabandi, K. Konolige, S. Levine, and V. Kumar. Deep dynamics models for learning
280 dexterous manipulation. In *Conference on Robot Learning*, pages 1101–1112. PMLR, 2020.
- 281 [11] T. Chen, J. Xu, and P. Agrawal. A system for general in-hand object re-orientation. In
282 *Conference on Robot Learning*, pages 297–307. PMLR, 2022.
- 283 [12] T. Chen, M. Tippur, S. Wu, V. Kumar, E. Adelson, and P. Agrawal. Visual dexterity: In-hand
284 dexterous manipulation from depth. *arXiv e-prints*, pages arXiv–2211, 2022.
- 285 [13] A. Handa, A. Allshire, V. Makoviychuk, A. Petrenko, R. Singh, J. Liu, D. Makoviichuk,
286 K. Van Wyk, A. Zhurkevich, B. Sundaralingam, et al. Dextreme: Transfer of agile in-hand
287 manipulation from simulation to reality. *arXiv preprint arXiv:2210.13702*, 2022.
- 288 [14] Z.-H. Yin, B. Huang, Y. Qin, Q. Chen, and X. Wang. Rotating without seeing: Towards in-hand
289 dexterity through touch. *arXiv preprint arXiv:2303.10880*, 2023.
- 290 [15] G. Khandate, M. Haas-Heger, and M. Ciocarlie. On the feasibility of learning finger-gaiting
291 in-hand manipulation with intrinsic sensing. In *2022 International Conference on Robotics and*
292 *Automation (ICRA)*, pages 2752–2758. IEEE, 2022.
- 293 [16] O. M. Andrychowicz, B. Baker, M. Chociej, R. Józefowicz, B. McGrew, J. Pachocki, A. Petron,
294 M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and
295 W. Zaremba. Learning dexterous in-hand manipulation. *The International Journal of Robotics*
296 *Research*, 39(1):3–20, 2020.
- 297 [17] H. Qi, A. Kumar, R. Calandra, Y. Ma, and J. Malik. In-hand object rotation via rapid motor
298 adaptation. In *Conference on Robot Learning*, pages 1722–1732. PMLR, 2023.

- 299 [18] A. Allshire, M. MittaI, V. Lodaya, V. Makoviychuk, D. Makoviichuk, F. Widmaier, M. Wüthrich,
300 S. Bauer, A. Handa, and A. Garg. Transferring dexterous manipulation from gpu simulation to
301 a remote real-world trifinger. In *2022 IEEE/RSJ International Conference on Intelligent Robots
302 and Systems (IROS)*, pages 11802–11809. IEEE, 2022.
- 303 [19] A. M. Okamura, N. Smaby, and M. R. Cutkosky. An overview of dexterous manipulation. In
304 *Proceedings 2000 ICRA. Millennium Conference. IEEE International Conference on Robotics
305 and Automation. Symposia Proceedings (Cat. No. 00CH37065)*, volume 1, pages 255–262.
306 IEEE, 2000.
- 307 [20] A. Bhatt, A. Sieler, S. Puhlmann, and O. Brock. Surprisingly robust in-hand manipulation: An
308 empirical study. *Robotics: Science and Systems (RSS)*, 2021.
- 309 [21] V. Kumar, A. Gupta, E. Todorov, and S. Levine. Learning dexterous manipulation policies from
310 experience and imitation. *arXiv preprint arXiv:1611.05095*, 2016.
- 311 [22] B. Calli, K. Srinivasan, A. Morgan, and A. M. Dollar. Learning modes of within-hand manip-
312 ulation. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages
313 3145–3151. IEEE, 2018.
- 314 [23] T. Ishihara, A. Namiki, M. Ishikawa, and M. Shimojo. Dynamic pen spinning using a high-speed
315 multifingered hand with high-speed tactile sensor. In *6th IEEE-RAS International Conference
316 on Humanoid Robots*, pages 258–263. IEEE, 2006.
- 317 [24] H. Van Hoof, T. Hermans, G. Neumann, and J. Peters. Learning robot in-hand manipulation
318 with tactile features. In *2015 IEEE-RAS 15th International Conference on Humanoid Robots
319 (Humanoids)*, pages 121–127. IEEE, 2015.
- 320 [25] B. Calli and A. M. Dollar. Vision-based model predictive control for within-hand precision
321 manipulation with underactuated grippers. In *2017 IEEE International Conference on Robotics
322 and Automation (ICRA)*, pages 2839–2845. IEEE, 2017.
- 323 [26] N. Furukawa, A. Namiki, S. Taku, and M. Ishikawa. Dynamic regrasping using a high-speed
324 multifingered hand and a high-speed vision system. In *Proceedings 2006 IEEE International
325 Conference on Robotics and Automation, 2006. ICRA 2006.*, pages 181–187. IEEE, 2006.
- 326 [27] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino,
327 M. Plappert, G. Powell, R. Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint
328 arXiv:1910.07113*, 2019.
- 329 [28] J. Tan, T. Zhang, E. Coumans, A. Iscen, Y. Bai, D. Hafner, S. Bohez, and V. Vanhoucke.
330 Sim-to-real: Learning agile locomotion for quadruped robots. *Robotics: Science and Systems
331 (RSS)*, 2018.
- 332 [29] X. Da, Z. Xie, D. Hoeller, B. Boots, A. Anandkumar, Y. Zhu, B. Babich, and A. Garg. Learning
333 a contact-adaptive controller for robust, efficient legged locomotion. In *Conference on Robot
334 Learning*, pages 883–894. PMLR, 2021.
- 335 [30] Z. Li, X. Cheng, X. B. Peng, P. Abbeel, S. Levine, G. Berseth, and K. Sreenath. Reinforce-
336 ment learning for robust parameterized locomotion control of bipedal robots. In *2021 IEEE
337 International Conference on Robotics and Automation (ICRA)*, pages 2811–2817. IEEE, 2021.
- 338 [31] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel. Asymmetric actor critic
339 for image-based robot learning. *arXiv preprint arXiv:1710.06542*, 2017.
- 340 [32] G. B. Margolis, G. Yang, K. Paigwar, T. Chen, and P. Agrawal. Rapid locomotion via reinforce-
341 ment learning. *Robotics: Science and Systems (RSS)*, 2022.
- 342 [33] G. B. Margolis, T. Chen, K. Paigwar, X. Fu, D. Kim, S. Kim, and P. Agrawal. Learning to jump
343 from pixels. In *Conference on Robot Learning*, pages 1025–1034. PMLR, 2022.

- 344 [34] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter. Learning quadrupedal locomotion
345 over challenging terrain. *Science robotics*, 5(47):eabc5986, 2020.
- 346 [35] A. Kumar, Z. Fu, D. Pathak, and J. Malik. Rma: Rapid motor adaptation for legged robots.
347 *Robotics: Science and Systems (RSS)*, 2021.
- 348 [36] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield. Deep object pose
349 estimation for semantic robotic grasping of household objects. *arXiv preprint arXiv:1809.10790*,
350 2018.
- 351 [37] D. Chen, B. Zhou, V. Koltun, and P. Krähenbühl. Learning by cheating. In *Conference on Robot
352 Learning*, pages 66–75. PMLR, 2020.
- 353 [38] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction
354 to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial
355 intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings,
356 2011.
- 357 [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
358 I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*,
359 30, 2017.
- 360 [40] O. Khatib. A unified approach for motion and force control of robot manipulators: The
361 operational space formulation. *IEEE Journal on Robotics and Automation*, 3(1):43–53, Feb.
362 1987. ISSN 2374-8710. doi:10.1109/JRA.1987.1087068. Conference Name: IEEE Journal on
363 Robotics and Automation.
- 364 [41] R. Tedrake and the Drake Development Team. Drake: Model-based design and verification for
365 robotics, 2019. URL <https://drake.mit.edu>.