

MimicTouch: Leveraging Multi-modal Human Tactile Demonstrations for Contact-rich Manipulation

Anonymous Author(s)

Affiliation

Address

email

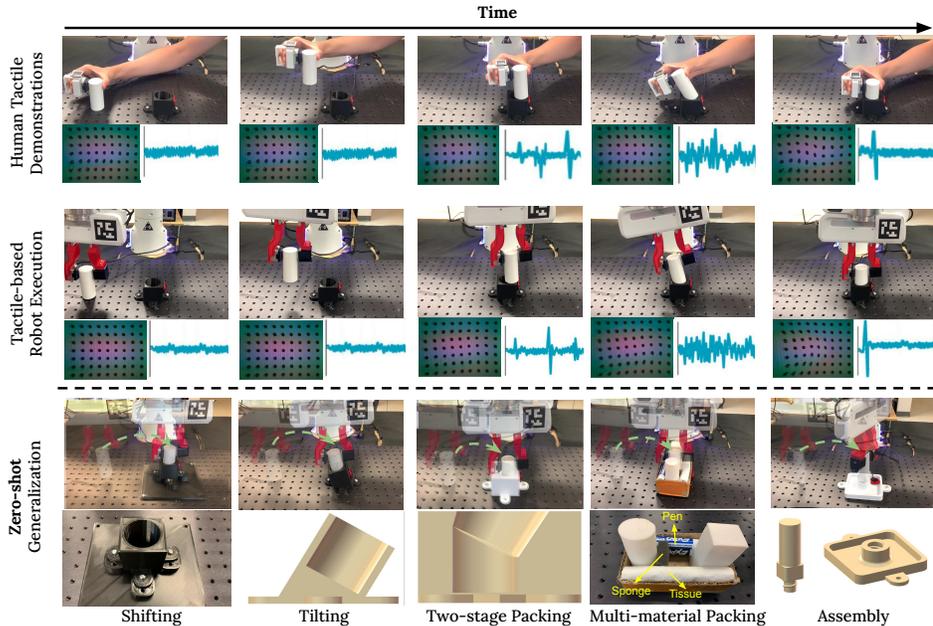


Figure 1: The first row shows the human tactile demonstrations, including the tactile and proprioception data. The second row shows the robot execution with tactile feedback. The third row below the dashed line describes the policy’s **zero-shot** generalization capability in five different domains.

1 **Abstract:** Tactile sensing is critical to fine-grained, contact-rich manipulation
2 tasks, such as insertion and assembly. Prior research has shown the possibility
3 of learning tactile-guided policy from teleoperated demonstration data. However,
4 to provide the demonstration, human users often rely on visual feedback to control
5 the robot. This creates a gap between the sensing modality used for controlling
6 the robot (visual) and the modality of interest (tactile). To bridge this gap,
7 we introduce “MimicTouch”, a novel framework for learning policies directly
8 from demonstrations provided by human users with their hands. The key
9 innovations are i) a human tactile data collection system which collects multi-
10 modal tactile dataset for learning human’s tactile-guided control strategy, ii) an
11 imitation learning-based framework for learning human’s tactile-guided control
12 strategy through such data, and iii) an online residual RL framework to bridge the
13 embodiment gap between the human hand and the robot gripper. Through comprehensive
14 experiments, we highlight the efficacy of utilizing human’s tactile-guided control
15 strategy to resolve contact-rich manipulation tasks. The project website is
16 at <https://sites.google.com/view/MimicTouch>.

17 **Keywords:** Tactile Sensing, Learning from Human, Imitation Learning

18 1 Introduction

19 Enabling robots to perform contact-rich tasks such as insertion remains a formidable challenge in
20 robotics. The primary reason is the complex dynamic interaction between the robot and the object

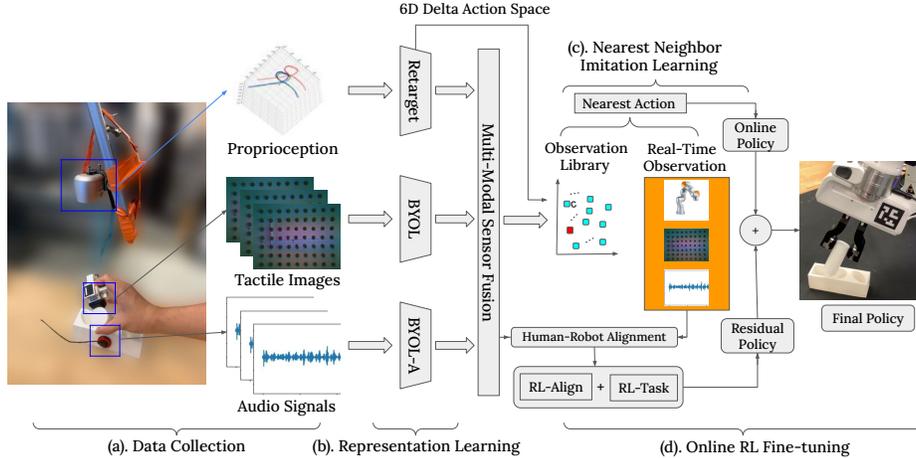


Figure 2: Illustration of the MimicTouch Framework. In part (a), we collect the multi-modal human tactile demonstrations. In part (b), we learn compact low-dimensional tactile representations. In part (c), we derive an offline policy through a non-parametric imitation learning method. In part (d), we refine the offline policy through online residual reinforcement learning on a physical robot.

21 which is influenced by various factors including intricate material properties and low tolerance for
 22 error. This necessitates an adaptive, data-centric insertion mechanism that utilizes real-time sensor
 23 feedback. Recent methods have heavily explored vision-based solutions to tackle this problem [1–
 24 4]. Notably, NVIDIA’s sim-to-real transfer approach [1] achieves success rates of up to 99.2%
 25 in transferring assembly tasks using their customized “Factory” simulator [2]. However, vision-
 26 based approaches can fall short when visual feedback is compromised by cluttered occlusions or
 27 bad lighting conditions.

28 Humans exhibit fine-grained manipulation skills through tactile sensing, which allows for successful
 29 insertions by solely using tactile feedback to generate complex, continuous, and precise motions [5].
 30 Motivated by this, recent studies collect demonstrations that combine various sensory inputs for pol-
 31 icy learning [6–8]. However, these methods assume a limited action space, e.g., only 3D translations,
 32 to compensate for the complexity of collecting dynamic demonstrations. They also heavily rely on
 33 robot teleoperation systems [9–11], which inevitably creates a gap between the visual sensing used
 34 for data collection and the recorded tactile sensing for policy learning. As a result, these methods
 35 are unable to emulate human’s tactile-guided control strategy for contact-rich tasks.

36 To tackle these challenges, we present “MimicTouch” (shown in Fig. 2), a novel framework that
 37 enables robots to learn human’s tactile-guided control strategy. Specifically, MimicTouch first intro-
 38 duces a human tactile-guided data collection system to gather multi-modal tactile feedback (tactile +
 39 audio) directly from human demonstrators. Next, it incorporates a representation learning model to
 40 capture task-specific sensor input features. These compact representations enhance the performance
 41 of subsequent imitation learning by abstracting essential sensory information. Then, it employs a
 42 non-parametric imitation learning method [12] to derive an offline policy from the collected human
 43 tactile demonstrations. Finally, it leverages online residual reinforcement learning (RL) to fine-tune
 44 the offline policy on the physical robot, aiming to bridge the embodiment gap between the human
 45 hand and the robot gripper and enrich contact reasoning.

46 We conduct comprehensive experiments on contact-rich insertion tasks to evaluate the offline pol-
 47 icy derived from demonstrations and the final policy fine-tuned via RL. We find that MimicTouch
 48 can collect tactile demonstrations more efficiently than teleoperation. More importantly, the Mim-
 49 icTouch policy can be effectively learned from such demonstrations and significantly outperforms
 50 the policy learned from teleoperation data. Additionally, we set up seven generalization tasks in five
 51 different domains and show the final policy exhibits superior **zero-shot** generalization capability.

52 2 Related Works

53 **Multi-modal tactile sensing.** Vision-based tactile sensing is integral to robotic manipulation as
 54 it excels at estimating local contact geometry and frictional properties [13–16]. It further enables
 55 imitation learning-based methods through policy observations [6, 7, 17] or reinforcement learning
 56 methods through reward signals [18–20] for various complex manipulation tasks. Additionally, it is

57 also widely used in shape reconstruction [21–23] and grasping [24–27]. On the other hand, audio-
58 based tactile sensors, such as contact microphones have also been demonstrated effective in robotics
59 applications such as manipulation [28], classification [29, 30], and dynamics modelling [31, 32].
60 These sensors can emulate nerve endings within human skin to better detect vibrations during tactile
61 interactions. Therefore, incorporating both sensor modalities can yield tactile feedback more akin
62 to human sensations, enabling the robot to learn a human-like tactile-guided control strategy.

63 **Learning from human demonstrations.** Learning from human demonstrations is a long-standing
64 research topic. One group of methods learns the robot behaviors from human videos [33–35]. How-
65 ever, these methods adopt only visual sensing, which can be easily collected by low-cost cameras
66 or accessed via online videos. As a result, they mostly focus on high-level scene reasoning rather
67 than fine-grained, contact-rich tasks, which often require tactile feedback for reliable execution. To
68 incorporate the tactile data into the demonstrations, recent works have turned to using robot teleop-
69 eration systems [6–8, 36–38], where a robot is equipped with all necessary sensors, including tactile
70 sensors, and is directly controlled by a human operator during task execution. This multi-modal
71 dataset will then be used to train robot policies. However, human operators must guide the robot
72 using visual feedback, thereby creating a gap between the visual sensing used for data collection and
73 the tactile sensing recorded for policy learning. Furthermore, collecting 6D dynamic motions for
74 contact-rich manipulations via teleoperation is challenging. Therefore, in this work, we propose to
75 collect human tactile demonstrations, in which the sensing gap is addressed and the demonstration
76 motions are more versatile and dynamic.

77 **Imitation learning.** Offline imitation learning (IL) is an effective strategy to learn robot policies
78 in the real world. We consider two classes of IL methods: parametric methods [39–41] and non-
79 parametric methods [12, 42, 43]. Parametric methods typically train neural networks to map ob-
80 servations to expert actions. While general in principle, they are prone to covariant shift and com-
81 pounding errors [44]. Our method instead adopts a non-parametric imitation learning method. These
82 methods constrain robot behaviors to the demonstrated data via techniques such as nearest-neighbor
83 lookup [12]. While they may be less general, they offer a safer alternative to their parametric coun-
84 terparts, which is crucial for the real-world contact-rich manipulation tasks considered in this work.

85 3 MimicTouch Framework

86 We aim to enable the robot to resolve contact-rich manipulation tasks by learning control strat-
87 egy from human tactile demonstrations. To achieve this, we propose a novel learning framework
88 named “MimicTouch”. It first introduces a human tactile-guided data collection system (Sec. 3.1)
89 to collect a multi-modal tactile dataset from human demonstrators. Then, to emulate the human’s
90 tactile-guided control strategy for successful robot execution, MimicTouch has three distinct learn-
91 ing phases. Firstly, it learns lower dimensional tactile representations from the human tactile demon-
92 strations in a self-supervised manner (Sec. 3.2). Next, it derives an offline policy with the learned
93 representations using a non-parametric imitation learning method [12] (Sec. 3.3). Lastly, it refines
94 the offline policy through online residual reinforcement learning (Sec. 3.4). Note that this refinement
95 phase is efficient and reliable as the offline policy encodes human’s tactile-guided control strategy.
96 The overall MimicTouch framework is shown in Fig. 2.

97 3.1 Collecting Human Tactile Demonstrations

98 To collect tactile demonstrations, current teleoperation systems have three key limitations: i). limited
99 scalability due to the need for a robot to collect demonstrations [10], ii). long training time and
100 expertise to become proficient with the system for fine-grained manipulation, and iii). sensing gap
101 between the visual sensing used for collection and recorded tactile sensing. To address these, our
102 key innovation is a system that enables humans to provide tactile demonstrations with their hands.
103 The system is elaborated in Fig. 6 (Appendix. A), and it collects the pose of human fingertips, tactile
104 images, and audio signals when human demonstrators perform contact-rich insertion tasks.

105 The data collection system consists of the following components. We use the RealSense camera
106 with Aruco Marker [45] for human fingertip pose tracking. The tracked fingertip poses are then
107 treated as the robot end-effector’s poses after calibration and filtering (Appendix. C.2). We also
108 use the GelSight Mini [46], a compact vision-based tactile sensor that is conveniently mounted onto
109 human fingertips using a custom fixture, to estimate the contacts between the object and the fingertip.

110 Notably, we only use one tactile sensor in our experiment setup instead of two. The Audio data,
111 which is helpful for manipulation tasks due to its sensitivity to contact vibration signals [28, 47],
112 is captured using the HOYUJI TD-11 piezo-electric contact microphone. Considering the potential
113 discrepancies in the mechanical vibrations between the human and the robot, the microphone is
114 placed at the base of the insertion hole to ensure signal consistency.

115 3.2 Learning Tactile Representation

116 The policy learning on high-dimensional sensor inputs struggles with real-world deployment due
117 to computational burden and sensor noise. Additionally, in this work, variations may exist be-
118 tween sensor inputs from human tactile demonstrations and real-time robot sensor feedback due to
119 discrepancies in finger-object contact force. Therefore, inspired by recent works that learn lower-
120 dimensional embeddings for image-guided imitation learning [6, 12, 48] which can discover the
121 appropriate features that are helpful for policy learning, we first learn the compact representation
122 for both tactile and audio sensor data using self-supervised learning methods (part (b) in Fig. 2).
123 Intuitively, it identifies a low-dimensional embedding space where differently augmented images,
124 such as tactile images or audio spectrum, are projected to a similar embedding. As a result, these
125 embeddings are more efficient for online processing and more robust to task-irrelevant sensor noise.
126 This learning phase consists of the following two parts.

127 **Data collection.** We collect task-specific tactile-audio data from the human demonstrator. The
128 dataset encompasses successful, failed, and sub-optimal demonstrations. For each, we segment the
129 audio data at 2Hz. In total, we collect 7657 tactile images and 1,000 audio segments. More details
130 are shown in Appendix. B

131 **Self-supervised learning.** We employ the Bootstrap Your Own Latent (BYOL) [49] for tactile
132 images and BYOL for audio (BYOL-A) for audio segments [50], since they have demonstrated
133 desired performance in computer vision [49], audio representation [50], and robotics [7, 28, 36]
134 tasks. Details about BYOL and BYOL-A are included in the Appendix. B.

135 3.3 Learning Offline Policy from Human Tactile Demonstrations

136 Leveraging the learned representations, we then learn the robot policy from the human tactile demon-
137 strations. Here, one unique challenge is that the human hand moves much faster than the robot,
138 resulting in sparse temporal observation-action samples (i.e., large action values per observation).
139 Also, human tactile demonstrations might partially be out-of-domain (OOD) demonstrations for
140 robot policy due to the embodiment gap (e.g., different motion capability and finger-object contact
141 forces). Therefore, executing a parametric policy might exhibit unreasonable robot behaviors as they
142 are prone to covariant shift [44] (see Sec. 4.2.1 for experimental validation). As a result, we use a
143 non-parametric imitation learning method [12] to ensure the execution efficacy of the policy learned
144 from human tactile demonstrations (part (c) in Fig. 2). In addition to the learning algorithm, data
145 pre-processing is necessary for synchronization and the details are included in Appendix. C.

146 **Non-parametric imitation learning.** We build our algorithm on the VINN framework [12] by
147 extending it to tactile-audio representation. At the i -th time step, the observations and actions are
148 denoted as $(o_i^T, o_i^A, o_i^{EE}, a_i)$, where o^T is the tactile representation, o^A is the audio representation,
149 o^{EE} is the robot end-effector pose, and the action a is defined by the $6D$ delta pose of the robot
150 end-effector, including a delta position and a delta Euler angle. Then, we extract tactile and audio
151 features (y_i^T, y_i^A) from (o_i^T, o_i^A) using the pre-trained representation encoders, respectively. These
152 tactile embeddings and the robot end-effector pose (y_i^T, y_i^A, o_i^{EE}) are formulated as the key features
153 of the demonstration library, with each associated with a corresponding action value a_i . Given the
154 varying scales of these inputs, we normalize them such that the maximum distance for each input
155 is unity in the library. In robot execution, for a given real-time observation $(\hat{o}_i^T, \hat{o}_i^A, \hat{o}_i^{EE})$, we first
156 obtain the query feature $(\hat{y}_i^T, \hat{y}_i^A, \hat{o}_i^{EE})$, and then search the demonstration library for a nearest-
157 neighbor-based action prediction.

158 3.4 Learning Residual Policy through Online Reinforcement Learning

159 The offline policy learned from human tactile demonstrations might not guarantee task success when
160 deployed on the physical robot. This could be due to: i). morphological differences between the

161 human hand and the robot gripper, ii). inaccurate fingertip tracking caused by fast movements, and
 162 iii). underexplored contact effects. Therefore, motivated by recent works using pure RL [17–20, 51]
 163 to learn tactile policies, we further leverage online reinforcement learning that allows in-domain
 164 robot interactions (part (d) of Fig. 2). It is noteworthy that the previous pure RL methods often
 165 generate quasi-static motions and utilize a limited action space [17–19, 51] because they learn from
 166 scratch without effective priors. On the contrary, we intend to leverage the best of both advantages
 167 by RL fine-tuning the offline policy learned from human tactile demonstrations.

168 Since it is infeasible to directly fine-tune the non-parametric policy, we instead learn a residual
 169 policy. Same as the offline policy, the input to the residual policy π_r is $(\hat{y}_t^T, \hat{y}_t^A, \delta_t^{EE})$ and output
 170 is the residual action a_i^r . Considering we use $6D$ continuous action space and sparse observation-
 171 action pairs (around 70 actions per trajectory), we opt for SAC [52] to handle the continuous action
 172 space with entropy regularization and to generate a replay buffer to increase the size of training data.
 173 The pseudocode of the RL training is included in Appendix. D. Finally, the robot action is the sum
 174 of the action generated from offline policy π_i and the action from the residual policy π_r .

175 Another critical component of residual RL is the reward design, which must balance exploitation and
 176 exploration. To address this, we combine an expert-aligned reward with a task-specific reward. The
 177 expert-aligned reward encourages a policy distribution that mimics the demonstrations, whereas the
 178 task-specific reward drives exploration to optimize the in-domain robot policy. More details about
 179 pseudocode, policy design, and rewards are included in Appendix. D.

180 4 Experiments

181 In this section, we first describe the experiment setting and the data collection throughput to high-
 182 light that MimicTouch can efficiently collect useful demonstrations (Sec. 4.1). Then we introduce
 183 the Offline Policy Evaluation to validate the efficacy of the offline policy and highlight the benefits
 184 of using human tactile demonstrations (Sec. 4.2), and the Online Policy Improvement and General-
 185 ization Evaluation to demonstrate the efficiency of learning the residual policy through online RL
 186 and the superior zero-shot generalization capability (Sec. 4.3).

187 **Hardware setting.** All experiments are conducted on a Franka Emika Panda Arm. For each task,
 188 the learned policy generates the 6-DoF pose command and then maps it to 7-DoF joint torque actions
 189 using an inverse kinematics solver and a low-level built-in controller.

190 **Teleoperation setting.** We compare our human tactile-guided data collection system with
 191 Spacemouse-based teleoperation, a popular teleoperation interface for manipulation tasks [9–11].
 192 To collect a similar number of observation-action pairs for each trajectory, we collect one robot
 193 state, one tactile image, and 0.5s audio segment at 5Hz. Since collecting teleoperation data necessi-
 194 tates considerable expertise, we allocate approximately 5 hours to practice with this system.

195 **Tasks.** We focus on two-piece insertion tasks that exemplify the challenge of many contact-rich
 196 manipulation settings. We 3D-print a cylinder and an insertion hole base and set up the same task
 197 environments for both data collection settings. An example of this task has been shown in Fig. 1

198 4.1 Human Tactile Demonstration Collection System

199 In this section, we demonstrate that Human Tactile Demonstrations can greatly improve data col-
 200 lection throughput for contact-rich manipulation tasks. We begin by determining the **usability** of
 201 a demonstration trajectory based on the following two metrics: i) the robot successfully inserts the
 202 object into the hole without any slipping or falling, and ii) the robot completes the task within 100
 203 actions. Using these criteria, we record the time length of collecting 20 usable demonstration tra-
 204 jectories by using our customized system (3.1) and teleoperation system (4). Then, we evaluate the
 205 data collection throughput (see Table. 1) in two metrics: i). the number of usable demonstrations
 206 collected per hour, and ii). the success rate for collecting usable demonstrations.

207 The results in Table. 1 support our insights: i). human tactile demonstrations can be collected sig-
 208 nificantly more efficiently than teleoperation systems for contact-rich tasks, and ii). human tactile
 209 demonstrations can seamlessly integrate human’s tactile feedback and motion capability, whereas
 210 teleoperation systems struggle to capture such dynamic tactile-guided motions. These factors to-
 211 gether result in much lower data collection efficiency and success rates of the teleoperation system.

Methods	Frequency	Success Rate
Human Tactile Demonstrations	104 traj/hr	83.3% (20/24)
Teleoperation	19 traj/hr	38.5% (20/52)

Table 1: Data collection throughput for human tactile demonstrations and teleoperation.

212 4.2 Offline Policy Evaluation

213 4.2.1 MimicTouch effectively learns from human tactile demonstrations

214 In this subsection, we evaluate the offline policy learned from the human tactile demonstrations
 215 using both VINN and a parametric imitation learning method. Specifically, we aim to test whether
 216 the offline policies can be deployed in real-world environments within the desired error tolerance,
 217 which is crucial for physical robot execution. To evaluate it, we used the data collection system (Sec.
 218 3.1) to collect 20 noise-free data sequences as the datasets to learn both offline policies. For testing,
 219 we gather another 5 data sequences with random noise to emulate the real-world environments
 220 (details are in Appendix. E). We then compute the mean square error loss (MSE Loss, defined
 221 in Appendix. F) to evaluate the policies on the testing set.

222 For the baseline parametric imitation learning method, we select the MULSA [6], which has been
 223 previously demonstrated effective in multisensory robot learning for insertion tasks. In our setting,
 224 we use the same sensor input (y_i^T, y_i^A, o_i^{EE}) as in the VINN method to generate the continuous $6D$
 225 delta action a_i . The same MSE loss is used for policy training and validation.

226 For both offline policies, we calculate the MSE losses between the generated action sequence and
 227 the ground truth action sequence in the testing set. We observed that the MSE loss from the VINN
 228 policy is **0.21**, which is significantly lower than that of MULSA policy (**1.53**). This suggests that
 229 the VINN policy can generate more accurate $6D$ continuous actions, indicating it is more suitable
 230 for subsequent online real-world RL fine-tuning. Notably, since the MSE Loss of MULSA is signif-
 231 icantly higher, we do not use it for further real-world experiments.

232 Additionally, we conduct the ablation study on the choice of multi-modal sensor inputs in Ap-
 233 pendix. G. The results suggest that multi-modal tactile feedback is crucial for the success of contact-
 234 rich insertion tasks, particularly during the insertion phase when most contact occurs.

235 4.2.2 Human tactile demonstrations trains better policies than teleoperated demonstrations

236 In this subsection, we compare the performance of the offline policies trained from human tactile
 237 demonstrations and teleoperation demonstrations. We collect 20 trajectories for each, and then we
 238 use the same VINN method to learn the offline policies for both sets of demonstrations. We evaluate
 239 the policies in two manners: i) the task success rate over 25 policy rollouts, and ii) the action serial
 240 numbers (right part of Fig. 3), i.e., the indexed numbers of the selected actions in the corresponding
 241 trajectory of the demonstration library, for each action of the rollout trajectories.

242 We first report that the task success rate of the offline policy rollouts from human tactile demonstra-
 243 tion is **40%**, whereas that of the policy rollouts from teleoperation data is only **12%**. Then, in the
 244 right part of Fig. 3, we show the mean and variance of the action serial numbers of three successful
 245 rollout trajectories for each offline policy. For the human tactile demonstrations, we observe a linear
 246 relationship with minimal variance. On the contrary, the action serial numbers for the teleoperation
 247 policy exhibit a non-linear relationship with a significantly larger variance, particularly during the
 248 insertion phase. This performance discrepancy arises because the majority of contacts occur dur-
 249 ing this phase, and the teleoperation lack of human tactile feedback is not well-suited for capturing
 250 these contact-rich events. Therefore, the indexes of the selected key features tend to be more disor-
 251 dered. A similar conclusion can be drawn from the screenshots (left part of Fig. 3) for one successful
 252 rollout of both policies. The Human Tactile Demonstrations policy exhibits dynamic tilting for ob-
 253 ject insertion, which emulates the demonstrated human’s tactile-guided control strategy (see Fig. 1).
 254 However, the teleoperation policy forcefully inserts the object with little orientation, indicating that
 255 teleoperation demonstrations capture less versatile motions that are necessary for contact-rich tasks.

256 Therefore, combined with the results in Sec. 4.1, MimicTouch not only efficiently collects human
 257 tactile demonstrations, but also enables effective policy learning using these demonstrations.

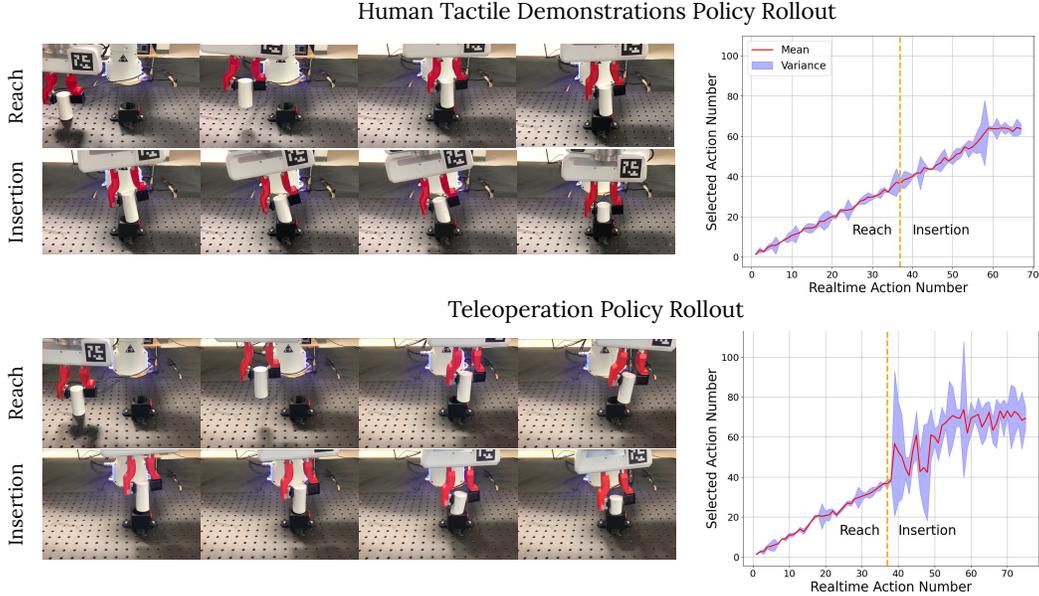


Figure 3: **Left:** Qualitative results for Human Tactile Demonstrations policy and Teleoperation policy. **Right:** Visualization of the action serial numbers for three trajectories generated by both policies. Solid red lines indicate mean trends and shaded areas show \pm standard deviations. The left side of the dashed orange line is the Reach phase, and the right side is the Insertion phase.

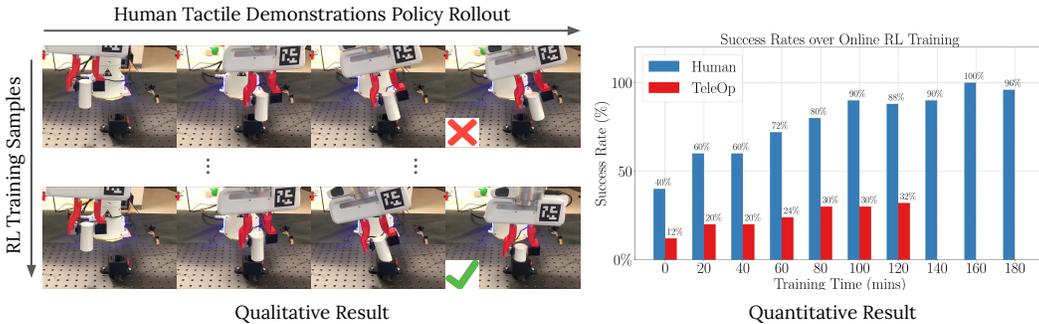


Figure 4: **Left:** Demonstrations of the online RL fine-tuning process, which further improves the task performance. **Right:** Quantitative task evaluations for offline policies learned from human tactile demonstrations (Human) and teleoperation (TeleOp) during online RL fine-tuning show that Human significantly outperforms TeleOp in terms of task success rate and RL training efficiency.

258 **4.3 Online Policy Improvement and Generalization to New Settings**

259 In this section, we evaluate the final policy trained through online reinforcement learning (RL). To
 260 ensure the robustness of the RL policy, we perform domain randomization on the robot’s starting
 261 position so that the initial object-hole contact is located differently. For the policy update, at each
 262 iteration, we use five newly collected trajectories along with another five randomly selected trajec-
 263 tories from the replay buffer.

264 **Online RL fine-tuning significantly and efficiently improves task performance.** We evaluate
 265 the trained policy every 20 minutes, approximately after every 13 RL epochs. After each hour,
 266 we compute the task success rates over 25 policy rollouts, in alignment with the offline policy
 267 evaluation. For other time instances, we only compute the task success rates over 10 policy rollouts
 268 to minimize sensor wear. The evaluation results are shown in Fig. 4, and we can observe that the
 269 policy can reach **96%** task success rate in **3 hours**. Additionally, the offline policy can reach 88%
 270 task success rate after 2-hour RL fine-tuning, which is significantly more training efficient than the
 271 policy learned from teleoperation, which could only achieve 32% task success rate at that time. This
 272 result supports the effectiveness of online RL fine-tuning, as it allows the robot to further interact



Figure 5: Setup of zero-shot generalization tasks.

with the task environment. Moreover, it once again highlights the importance of using human tactile demonstrations since the offline policy learned from teleoperation demonstration exhibit significant training inefficiency in the subsequent online RL fine-tuning.

MimicTouch policy exhibits superior zero-shot generalization capability. We evaluate the zero-shot generalizability of the MimicTouch final policy. We consider the following generalization settings: i). *Shifting Positions*: an insertion task with the hole positions shifted for 0.8cm in either $\pm x$ or $\pm y$ horizontal directions, ii). *Tilting Angles*: an insertion task with the hole angle tilted for 10 degrees or 20 degrees, iii). *Two-stage Dense Packing*: a two-stage dense insertion task, which requires the robot to perform two consecutive insertion alignments, iv). *Multi-material Dense Packing*: an insertion task where the hole contains multiple objects, such as a pen, tissues, or a sponge, and v) *Furniture Assembly* [53, 54]: an insertion task, which requires the robot to insert and adjust the leg into a small hole in a table for screwing. Each setting is depicted in Fig. 5. Additionally, to demonstrate the complexity of these tasks, we introduce another baseline: *Openloop Policy*, where we collect five successful insertion trajectories in the initial setting and execute each of these trajectories for those generalization tasks five times for each of the tasks. See Appendix. H for the details of each task and the policy evaluation process.

Policy	Shift	10°	20°	Two-stage	Rigid	Soft	Assem (I)	Assem (A)	Assem
<i>Openloop</i>	60%	56%	40%	52%	52%	36%	32%	37.5%	12%
MimicTouch	92.5%	92%	80%	88%	80%	64%	76%	68.4%	52%

Table 2: Task success rate for each generalization task. Specifically, Shift is *Shifting Position*, 10° and 20° are different tilted angles for *Tilting Angles*, Two-stage is *two-stage Dense Packing*, Rigid means the object contacts the pen in *Multi-material Dense Packing*, Soft means the object contacts the tissue and sponge in *Multi-material Dense Packing*, Assem (I), Assem (A), and Assem are the insertion, adjustment, and overall results for *Furniture Assembly* respectively.

We report the task success rate of both policies in Table. 2. Based on these results, we can observe that the MimicTouch policy can significantly outperform all the openloop policies in all those generalization tasks. Also, since the *Openloop* policy has lower success rates in all the generalization tasks, it indicates the robustness of the MimicTouch policy in different challenging generalization domains. We also include the qualitative results and detailed analysis in Appendix. I.

5 Limitation and Future Work

MimicTouch pioneers the pathway to learning human tactile-guided control strategies from human tactile demonstrations. However, it still has several limitations for future improvements. Firstly, MimicTouch still requires several hours to RL fine-tune the policy to address the embodiment gap and enrich the contact reasoning. We will explore a better representation learning method to further reduce the gap between humans and robots. Secondly, MimicTouch is task-specific and can not directly generalize human’s tactile-guided control strategy to other contact-rich manipulation tasks. One potential solution is to learn a generalizable tactile-based dynamic model, rather than a task-specific control policy. Finally, the method of learning to perform contact-rich manipulation tasks from human tactile demonstrations could be extended to other robot tasks, such as dexterous manipulation, bimanual manipulation, and deformable object manipulation.

6 Conclusion

We presented MimicTouch, a multi-modal imitation learning framework that: i) enables humans to perform tactile demonstrations with their hands without a robot in the loop ii) learns from such demonstrations and safety transfer to robot with non-parametric imitation learning, and iii) improves the policy performance with residual-based online RL to bridge the human-robot embodiment gap. We show that MimicTouch enables high-throughput data collection and achieves high success rate and generalization across a wide range of two-piece insertion and assembly tasks.

References

- [1] B. Tang, M. A. Lin, I. Akinola, A. Handa, G. S. Sukhatme, F. Ramos, D. Fox, and Y. Narang. Industreal: Transferring contact-rich assembly tasks from simulation to reality. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- [2] Y. Narang, K. Storey, I. Akinola, M. Macklin, P. Reist, L. Wawrzyniak, Y. Guo, A. Moravanzky, G. State, M. Lu, A. Handa, and D. Fox. Factory: Fast contact for robotic assembly. In *Proceedings of Robotics: Science and Systems (RSS)*, 2022.
- [3] G. Schoettler, A. Nair, J. Luo, S. Bahl, J. Aparicio Ojea, E. Solowjow, and S. Levine. Deep reinforcement learning for industrial insertion tasks with visual inputs and natural rewards. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5548–5555, 2020. doi:10.1109/IROS45743.2020.9341714.
- [4] G. Schoettler, A. Nair, J. A. Ojea, S. Levine, and E. Solowjow. Meta-reinforcement learning for robotic industrial insertion tasks. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 9728–9735, 2020. doi:10.1109/IROS45743.2020.9340848.
- [5] G. Lengyel, G. Žalalytė, A. Pantelides, J. N. Ingram, J. Fiser, M. Lengyel, and D. M. Wolpert. Unimodal statistical learning produces multimodal object-like representations. *eLife*, 8:e43942, may 2019. ISSN 2050-084X. doi:10.7554/eLife.43942. URL <https://doi.org/10.7554/eLife.43942>.
- [6] H. Li, Y. Zhang, J. Zhu, S. Wang, M. A. Lee, H. Xu, E. Adelson, L. Fei-Fei, R. Gao, and J. Wu. See, hear, and feel: Smart sensory fusion for robotic manipulation. In *CoRL*, 2022.
- [7] I. Guzey, B. Evans, S. Chintala, and L. Pinto. Dexterity from touch: Self-supervised pre-training of tactile representations with robotic play. In *7th Annual Conference on Robot Learning*, 2023. URL <https://openreview.net/forum?id=EXQ0eXtX30W>.
- [8] A. Mandlekar, Y. Zhu, A. Garg, J. Booher, M. Spero, A. Tung, J. Gao, J. Emmons, A. Gupta, E. Orbay, S. Savarese, and L. Fei-Fei. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In A. Billard, A. Dragan, J. Peters, and J. Morimoto, editors, *Proceedings of The 2nd Conference on Robot Learning*, volume 87 of *Proceedings of Machine Learning Research*, pages 879–893. PMLR, 29–31 Oct 2018. URL <https://proceedings.mlr.press/v87/mandlekar18a.html>.
- [9] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- [10] C. Chi, Z. Xu, C. Pan, E. Cousineau, B. Burchfiel, S. Feng, R. Tedrake, and S. Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. In *Proceedings of Robotics: Science and Systems (RSS)*, 2024.
- [11] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. Viola: Imitation learning for vision-based manipulation with object proposal priors. *6th Annual Conference on Robot Learning (CoRL)*, 2022.
- [12] J. Pari, N. M. Shafiullah, S. P. Arunachalam, and L. Pinto. The surprising effectiveness of representation learning for visual imitation. In *Proceedings of Robotics: Science and Systems (RSS)*, 2022.
- [13] W. Yuan, S. Dong, and E. H. Adelson. Gelsight: High-resolution robot tactile sensors for estimating geometry and force. *Sensors*, 17(12), 2017. ISSN 1424-8220. doi:10.3390/s17122762. URL <https://www.mdpi.com/1424-8220/17/12/2762>.
- [14] E. Donlon, S. Dong, M. Liu, J. Li, E. Adelson, and A. Rodriguez. Gelslim: A high-resolution, compact, robust, and calibrated tactile-sensing finger. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1927–1934, 2018. doi:10.1109/IROS.2018.8593661.

- 359 [15] D. Ma, E. Donlon, S. Dong, and A. Rodriguez. Dense tactile force estimation using gelslim
360 and inverse fem. In *2019 International Conference on Robotics and Automation (ICRA)*, pages
361 5418–5424, 2019. doi:10.1109/ICRA.2019.8794113.
- 362 [16] S. Dong and A. Rodriguez. Tactile-based insertion for dense box-packing. In *2019 IEEE/RSJ*
363 *International Conference on Intelligent Robots and Systems (IROS)*, pages 7953–7960, 2019.
364 doi:10.1109/IROS40897.2019.8968204.
- 365 [17] M. A. Lee, Y. Zhu, K. Srinivasan, P. Shah, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg.
366 Making sense of vision and touch: Self-supervised learning of multimodal representations
367 for contact-rich tasks. In *2019 IEEE International Conference on Robotics and Automation*
368 *(ICRA)*, 2019. URL <https://arxiv.org/abs/1810.10191>.
- 369 [18] S. Dong, D. K. Jha, D. Romeres, S. Kim, D. Nikovski, and A. Rodriguez. Tactile-rl for inser-
370 tion: Generalization to objects of unknown geometry. In *2021 IEEE International Conference*
371 *on Robotics and Automation (ICRA)*, pages 6437–6443, 2021. doi:10.1109/ICRA48506.2021.
372 9561646.
- 373 [19] S. Kim and A. Rodriguez. Active extrinsic contact sensing: Application to general peg-in-hole
374 insertion. In *ICRA*, 2022.
- 375 [20] Z.-H. Yin, B. Huang, Y. Qin, Q. Chen, and X. Wang. Rotating without seeing: Towards in-hand
376 dexterity through touch. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- 377 [21] S. Suresh, Z. Si, J. G. Mangelson, W. Yuan, and M. Kaess. Shapemap 3-d: Efficient shape
378 mapping through dense touch and vision. In *2022 International Conference on Robotics and*
379 *Automation (ICRA)*, pages 7073–7080, 2022. doi:10.1109/ICRA46639.2022.9812040.
- 380 [22] E. J. Smith, R. Calandra, A. Romero, G. Gkioxari, D. Meger, J. Malik, and M. Drozdal. 3d
381 shape reconstruction from vision and touch. *arXiv preprint arXiv:2007.03778*, 2020.
- 382 [23] E. Smith, D. Meger, L. Pineda, R. Calandra, J. Malik, A. Romero Soriano, and
383 M. Drozdal. Active 3d shape reconstruction from vision and touch. In M. Ran-
384 zato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in*
385 *Neural Information Processing Systems*, volume 34, pages 16064–16078. Curran Asso-
386 ciates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2021/file/8635b5fd6bc675033fb72e8a3ccc10a0-Paper.pdf)
387 [2021/file/8635b5fd6bc675033fb72e8a3ccc10a0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/8635b5fd6bc675033fb72e8a3ccc10a0-Paper.pdf).
- 388 [24] R. Calandra, A. Owens, D. Jayaraman, J. Lin, W. Yuan, J. Malik, E. H. Adelson, and S. Levine.
389 More than a feeling: Learning to grasp and regrasp using vision and touch. *IEEE Robotics and*
390 *Automation Letters*, 3(4):3300–3307, oct 2018. doi:10.1109/lra.2018.2852779. URL <https://doi.org/10.1109/lra.2018.2852779>.
- 392 [25] L. Pinto, D. Gandhi, Y. Han, Y.-L. Park, and A. Gupta. The curious robot: Learning visual
393 representations via physical interactions. In *European Conference on Computer Vision*, pages
394 3–18. Springer, 2016.
- 395 [26] R. Calandra, A. Owens, M. Upadhyaya, W. Yuan, J. Lin, E. H. Adelson, and S. Levine. The
396 feeling of success: Does touch sensing help predict grasp outcomes? In S. Levine, V. Van-
397 houcke, and K. Goldberg, editors, *Proceedings of the 1st Annual Conference on Robot Learn-*
398 *ing*, volume 78 of *Proceedings of Machine Learning Research*, pages 314–323. PMLR, 13–15
399 Nov 2017. URL <https://proceedings.mlr.press/v78/calandra17a.html>.
- 400 [27] Y. Han, K. Yu, R. Batra, N. Boyd, C. Mehta, T. Zhao, Y. She, S. Hutchinson, and Y. Zhao.
401 Learning generalizable vision-tactile robotic grasping strategy for deformable objects via trans-
402 former. *arXiv preprint arXiv:2112.06374*, 2023.
- 403 [28] A. Thankaraj and L. Pinto. That sounds right: Auditory self-supervision for dynamic robot ma-
404 nipulation. In *7th Annual Conference on Robot Learning*, 2023. URL [https://openreview.](https://openreview.net/forum?id=sLhk0keeiseH)
405 [net/forum?id=sLhk0keeiseH](https://openreview.net/forum?id=sLhk0keeiseH).
- 406 [29] J. Sinapov, M. Wiemer, and A. Stoytchev. Interactive learning of the acoustic properties of
407 household objects. In *2009 IEEE International Conference on Robotics and Automation*, pages
408 2518–2524, 2009. doi:10.1109/ROBOT.2009.5152802.

- 409 [30] J. Christie and N. Kottege. Acoustics based terrain classification for legged robots. In *2016*
410 *IEEE International Conference on Robotics and Automation (ICRA)*, pages 3596–3603, 2016.
411 doi:10.1109/ICRA.2016.7487543.
- 412 [31] C. Matl, Y. Narang, D. Fox, R. Bajcsy, and F. Ramos. Stressd: Sim-to-real from sound for
413 stochastic dynamics. In J. Kober, F. Ramos, and C. Tomlin, editors, *Proceedings of the 2020*
414 *Conference on Robot Learning*, volume 155 of *Proceedings of Machine Learning Research*,
415 pages 935–958. PMLR, 16–18 Nov 2021. URL [https://proceedings.mlr.press/v155/
416 mat121a.html](https://proceedings.mlr.press/v155/mat121a.html).
- 417 [32] D. Gandhi, A. Gupta, and L. Pinto. Swoosh! rattle! thump! – actions that sound. In *Proceed-*
418 *ings of Robotics: Science and Systems (RSS)*, 2020.
- 419 [33] P. Sharma, D. Pathak, and A. Gupta. Third-person visual imitation learning via decoupled
420 hierarchical controller. *Advances in Neural Information Processing Systems*, 32, 2019.
- 421 [34] H. Bharadhwaj, A. Gupta, S. Tulsiani, and V. Kumar. Zero-shot robot manipulation from
422 passive human videos. *arXiv preprint arXiv:2302.02011*, 2023.
- 423 [35] H. Xiong, Q. Li, Y.-C. Chen, H. Bharadhwaj, S. Sinha, and A. Garg. Learning by watching:
424 Physical imitation of manipulation skills from human videos. In *2021 IEEE/RSJ International*
425 *Conference on Intelligent Robots and Systems (IROS)*, pages 7827–7834. IEEE, 2021.
- 426 [36] S. P. Arunachalam, I. Güzey, S. Chintala, and L. Pinto. Holo-dex: Teaching dexterity with
427 immersive mixed reality. In *2023 IEEE International Conference on Robotics and Automation*
428 *(ICRA)*, pages 5962–5969, 2023. doi:10.1109/ICRA48891.2023.10160547.
- 429 [37] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mimicplay:
430 Long-horizon imitation learning by watching human play. In *7th Annual Conference on Robot*
431 *Learning*, 2023. URL <https://openreview.net/forum?id=hRZ1YjDZmTo>.
- 432 [38] M. Xu, Z. Xu, C. Chi, M. Veloso, and S. Song. Xskill: Cross embodiment skill discovery. In
433 *Conference on Robot Learning*, 2023.
- 434 [39] J. Ho and S. Ermon. Generative adversarial imitation learning. In D. Lee, M. Sugiyama,
435 U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Sys-*
436 *tems*, volume 29. Curran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/
437 paper_files/paper/2016/file/cc7e2b878868cbae992d1fb743995d8f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/cc7e2b878868cbae992d1fb743995d8f-Paper.pdf).
- 438 [40] M. Xie, A. Handa, S. Tyree, D. Fox, H. Ravichandar, N. D. Ratliff, and K. V. Wyk. Neural ge-
439 ometric fabrics: Efficiently learning high-dimensional policies from demonstration. In K. Liu,
440 D. Kulic, and J. Ichnowski, editors, *Proceedings of The 6th Conference on Robot Learning*,
441 volume 205 of *Proceedings of Machine Learning Research*, pages 1355–1367. PMLR, 14–18
442 Dec 2023.
- 443 [41] Y. Han, M. Xie, Y. Zhao, and H. Ravichandar. On the utility of koopman operator theory in
444 learning dexterous manipulation skills. In *7th Annual Conference on Robot Learning*, 2023.
445 URL <https://openreview.net/forum?id=pw-OTIYrGa>.
- 446 [42] Y. Huang, L. Rozo, J. Silvério, and D. G. Caldwell. Non-parametric imitation learning of
447 robot motor skills. In *2019 International Conference on Robotics and Automation (ICRA)*,
448 pages 5266–5272, 2019. doi:10.1109/ICRA.2019.8794267.
- 449 [43] M. Vaandrager, R. Babuška, L. Buşoniu, and G. A. Lopes. Imitation learning with non-
450 parametric regression. In *Proceedings of 2012 IEEE International Conference on Automation,*
451 *Quality and Testing, Robotics*, pages 91–96, 2012. doi:10.1109/AQTR.2012.6237681.
- 452 [44] H. Ravichandar, A. S. Polydoros, S. Chernova, and A. Billard. Recent advances in robot
453 learning from demonstration. *Annual review of control, robotics, and autonomous systems*, 3:
454 297–330, 2020.
- 455 [45] S. Garrido-Jurado, R. Muñoz-Salinas, F. J. Madrid-Cuevas, and M. J. Marín-Jiménez. Auto-
456 matic generation and detection of highly reliable fiducial markers under occlusion. *Pattern*
457 *Recognit.*, 47:2280–2292, 2014. URL [https://api.semanticscholar.org/CorpusID:
458 12519167](https://api.semanticscholar.org/CorpusID:12519167).

- 459 [46] Gelsight mini. <https://www.gelsight.com/gelsightmini/>.
- 460 [47] M. Du, O. Y. Lee, S. Nair, and C. Finn. Play it by ear: Learning skills amidst occlusion through
461 audio-visual imitation learning. In *Proceedings of Robotics: Science and Systems (RSS)*, 2022.
- 462 [48] S. Haldar, J. Pari, A. Rai, and L. Pinto. Teach a robot to fish: Versatile imitation from one
463 minute of demonstrations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- 464 [49] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch,
465 B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, B. Piot, k. kavukcuoglu, R. Munos,
466 and M. Valko. Bootstrap your own latent - a new approach to self-supervised learning.
467 In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances
468 in Neural Information Processing Systems*, volume 33, pages 21271–21284. Curran As-
469 sociates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/
470 2020/file/f3ada80d5c4ee70142b17b8192b2958e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/f3ada80d5c4ee70142b17b8192b2958e-Paper.pdf).
- 471 [50] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino. Byol for audio: Self-
472 supervised learning for general-purpose audio representation. In *2021 International Joint
473 Conference on Neural Networks (IJCNN)*, pages 1–8, 2021. doi:10.1109/IJCNN52387.2021.
474 9534474.
- 475 [51] J. Ding, C. Wang, and C. Lu. Transferable force-torque dynamics model for peg-in-hole task,
476 2019.
- 477 [52] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy
478 deep reinforcement learning with a stochastic actor. In J. Dy and A. Krause, editors,
479 *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Pro-
480 ceedings of Machine Learning Research*, pages 1861–1870. PMLR, 10–15 Jul 2018. URL
481 <https://proceedings.mlr.press/v80/haarnoja18b.html>.
- 482 [53] M. Heo, Y. Lee, D. Lee, and J. J. Lim. Furniturebench: Reproducible real-world benchmark
483 for long-horizon complex manipulation. In *Robotics: Science and Systems*, 2023.
- 484 [54] H. Lin, R. Corcodel, and D. Zhao. Generalize by touching: Tactile ensemble skill transfer for
485 robotic furniture assembly, 2024.
- 486 [55] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016
487 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
488 doi:10.1109/CVPR.2016.90.

Appendices

490 A Data Collection

In this section, we show the novel data collection system in Fig. 6.

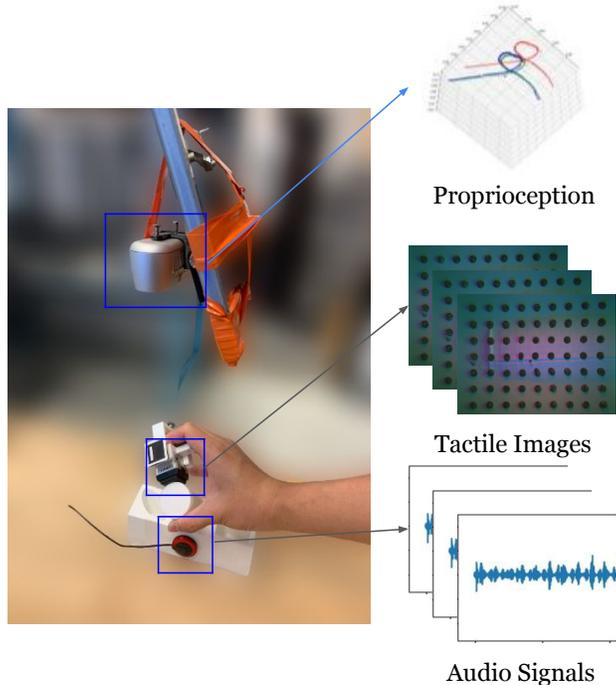


Figure 6: The human tactile demonstrations collection system.

491

492 B Representation Learning

493 **Data Collection** For Representation Learning, we collect a large task-specific dataset that contains
 494 7657 tactile images and 1000 audio sequences. We collect 100 trajectories, each of them approxi-
 495 mately five seconds long and containing around 70 tactile images and 10 audio sequences. These
 496 trajectories contain various data qualities, which include successful cases, failure cases, and sub-
 497 optimal cases. In detail, successful cases refer to the cases that human finishes the task with smooth
 498 trajectories; failure cases mean that human did not insert the object successfully or used more than
 499 five seconds to finish this task; sub-optimal cases mean that human used unnecessary motions to
 500 finish the insertion task.

501 **BYOL BYOL** [49] generates two augmented views, $v \triangleq t(x)$ and $v' \triangleq t'(x)$, from a given x by
 502 applying image augmentations $t \sim \mathcal{T}$ and $t' \sim \mathcal{T}'$ respectively, where \mathcal{T} and \mathcal{T}' represent distinct
 503 augmentation distributions. The architecture of BYOL comprises a primary encoder f_θ and a target
 504 encoder f_ξ , where the latter being an exponential moving average of the former. Given the aug-
 505 mented views v and v' , they are processed to yield representations y and y' . These representations
 506 are subsequently transformed by projectors g_θ and g_ξ to produce higher-dimensional vectors z and
 507 z' . The primary encoder and its associated projector are designed to predict the output from the tar-
 508 get projector, resulting in $q_\theta(z_\theta)$ and $sg(z'_\xi)$. The model’s output consists of l_2 -normalized versions
 509 of these predictions, which are trained using a similarity loss function. Post-training, the encoder f_θ
 510 is utilized for feature extraction from observations.

511 To utilize BYOL in tactile images, we scale the tactile image up to 256x256 to work with standard
 512 image encoders. We use the ResNet [55] architecture, also starting with pre-trained weights. Un-
 513 like SSL (self-supervised learning) techniques used in visual images, we only apply the Gaussian

514 blur and small center-resized crop augmentations, since other augmentations such as color jitter and
515 grayscale would violate the assumption that augmentations do not change the tactile signal signifi-
516 cantly. For each input, the trained model will generate a 1×2048 representation vector.

517 **Audio Representation Learning** BYOL-A [50] is an extended version of BYOL to audio rep-
518 resentation learning, processing log-scaled mel-spectrograms through a specialized augmentation
519 module. To utilize BYOL-A in our audio data, we down-sampled signals from 44.1kHz to 16kHz,
520 with a window size of 64 ms, a hop size of 10 ms, and mel-spaced frequency bins $F = 64$ in the range
521 60–7,800 Hz. Then, the Pre-Normalization step stabilizes the input audio for subsequent augmenta-
522 tions. Once normalized, the Mixup step introduces contrasts in the audio’s background, defined by
523 the log-mixup-exp formula:

$$\tilde{x}_i = \log((1 - \lambda) \exp(x_i) + \lambda \exp(x_k))$$

524 where x_k is a mixing counterpart and λ is a ratio from a uniform distribution. The next one is the
525 RRC block, an augmentation technique, that captures content details and emulates pitch shifts and
526 time stretches. For each input, the trained model will generate a 1×2048 representation vector.

527 C Data Pre-processing

528 C.1 Sensor Data Alignment

529 Each sensor operates at different frequency: i) RealSense operates at 60 Hz with a resolution of
530 320x240 pixels, ii) GelSight Mini streams tactile images at 15 Hz with 400x300 pixel resolution,
531 and iii) HOYUJI TD-11 piezo-electric contact microphone has a 44.1kHz sampling rate, and the
532 audio data is segmented at 2Hz.

533 Therefore, to ensure synchronization across our sensors, we first address the disparate sampling
534 rates of the fingertip poses, tactile images, and audio sequences, which are 60Hz, 15Hz, and 2Hz,
535 respectively. Specifically, we downsample the fingertip poses to 15 Hz. For the audio data, instead
536 of collecting entirely new 0.5-second segments, we record the extended audio signals at intervals
537 of every 0.07 seconds. As a result, the new 0.5s segment has a new-collected 0.07s interval and
538 an old overlapped 0.43s segment in the past, which results in an overlap of 0.43 seconds between
539 consecutive audio segments. Therefore, all sensor inputs are sampled at 15Hz.

540 C.2 Calibration and Filtering of Fingertip Poses

541 The 6D human fingertip poses extracted from the Aruco marker include 3D positions along with
542 rotation vectors. To use these fingertip poses as the end-effector’s poses for robot policy learning,
543 we need to address two problems: i). developing a calibration method to align the data collection
544 system with the robot execution system, and ii). implementing a filtering method to generate smooth
545 trajectories.

546 **Calibration** Given that data collection and robot experiments occur in disparate scenarios, it is
547 crucial to align our human-centric data collection system with the physical robot system. Initially,
548 we record the distance between the object (starting point) and the base (ending point) within the data
549 collection system and replicate this setup in the robot environment. Following this, six equidistant
550 positions between the starting and ending points are identified within both systems. The object is
551 gripped at these predetermined positions using both hands and the robot’s end-effector so that we
552 can capture the corresponding poses. In this calibration process, the hand poses, denoted as the
553 “Eye” in the calibration function, are referenced to the camera frame, while the end-effector poses,
554 represented as the “Hand” in the calibration system, are referenced to the robot frame. Conclusively,
555 we employ the `calibrateHandEye` function from OpenCV, using the six captured poses, to calibrate
556 these two frames (camera frame and robot frame).

557 **Filtering** Given the inherent noise and occasional outliers in the poses obtained from the RealSense
558 and Aruco markers, it is imperative to implement post-processing techniques to ensure the quality
559 and smoothness of the trajectories. For each pose sequence, outliers are detected by sorting the
560 values of each delta transformation (i.e., the delta translations and the delta Euler angles). The
561 Interquartile Range (IQR) method is employed to establish the upper and lower bounds, which are
562 then used to identify outliers. The IQR is defined as: $IQR = Q_3 - Q_1$ where Q_3 and Q_1 are the
563 third and first quartiles, respectively. Outliers are replaced using a median filter with a window size

Algorithm 1 Online Residual Reinforcement Learning

```
1: Input: offline policy  $\pi_i$ , randomly initialized residual policy  $\pi_r$ 
2: Input: step size sequences  $\{\beta_t\}$ , number of iterations  $K$ , Replay Buffer  $D$ 
3: Initialize replay buffer  $D$  with pre-collected data
4: for  $k = 1$  to  $K$  do
5:   Sample mini-batch  $D_k$  from Replay Buffer  $D$ 
6:   Obtain current trajectory  $C_k$  by executing  $\pi_i + \pi_r$ 
7:   To collect more data in  $D$ :  $D \leftarrow D \cup C_k$ 
8:   Combine  $D_k$  and  $C_k$  to form batch  $B_k$  for update
9:   for all  $(s, a_i, r, s') \in B_k$  do
10:      $a_i \leftarrow \pi_i(s)$  ▷ Obtain latent action
11:      $a_r \leftarrow \pi_r(s)$  ▷ Obtain residual action
12:      $\hat{a} \leftarrow a_i + a_r$  ▷ Combine latent and residual actions
13:      $Q_r \leftarrow Q_r(s, a_i) + r + \gamma Q_r(s', \pi_i(s'))$ 
14:      $\pi_r \leftarrow \pi_r - \alpha \nabla_{\pi_r} L(\pi_r)$  ▷ Update residual policy with gradient step
15:   end for
16: end for
17: Return: Trained residual policy  $\pi_r$ 
```

564 of 3. To enhance the temporal consistency of the estimated hand and object pose, a digital low-pass
565 filter is applied to eliminate high-frequency noise. Specifically, the filter has a sampling frequency
566 of 5Hz and a cutoff frequency of 2Hz. The low-pass filter can be represented as: $H(f) = \frac{1}{1+(\frac{f}{f_c})^2}$
567 where f is the sampling frequency and f_c is the cutoff frequency.

568 D Details of RL training

569 **Training pipeline** The overall pseudo-code for RL Training is given in Alg. 1.

570 **RL policy details** For the residual policy π_r within our framework, we employ the Soft Actor-
571 Critic (SAC) algorithm with an MLP architecture. The training strategy aims to effectively combine
572 reinforcement learning principles with residual corrections, thus enhancing the overall performance
573 of the system. The following formula represents the objective for training the residual policy:

$$\pi_r = \operatorname{argmax}_{\pi} \{ \mathbb{E}_{(s,a) \sim D} [Q(s, a + \pi(s)) - \alpha \log \pi(a|s)] \}$$

- 574 • $Q(s, a + \pi(s))$: The Q-value function, which estimates the value of executing the residual
575 action $\pi(s)$ in addition to the base action a in the state s . This represents the total action
576 influenced by both the offline policy and the residual corrections suggested by π_r .
- 577 • $\mathbb{E}_{(s,a) \sim D}$: The expectation over state-action pairs sampled from the replay buffer D , which
578 contains data from both past experiences and current new explorations.
- 579 • $\alpha \log \pi(a|s)$: The entropy regularization term for the policy π , which encourages explo-
580 ration by penalizing the certainty of the policy’s action selection. This term is crucial in
581 SAC to ensure sufficient exploration and avoid premature convergence to suboptimal poli-
582 cies.
- 583 • $\pi(a|s)$: The policy network (MLP) outputs the probability distribution over actions given
584 the state s , from which the action a is sampled.

585 This formula ensures that the residual policy π_r learns to adjust the actions generated by the offline
586 policy by optimizing the SAC objective. It balances the maximization of expected returns (via
587 Q-values) and the maintenance of behavioral diversity (via entropy regularization), allowing π_r to
588 adapt and refine actions based on real-time environmental feedback and historical data from the
589 replay buffer.

590 **RL Reward Design** We will give a detailed explanation for each component in our reward design.

591 **Distance Reward:**

$$d = 1 - \tanh(10.0 * ||distance||_2)$$

592 where the *distance* is between the current position of the gripper center and the target gripper center.

593 **Orientation Reward:**

$$o = 1 - \tanh(7.5 * ||diff_ori||_2)$$

594 where the *diff_ori* stands as the quaternion difference between the current gripper orientation and
595 the target gripper orientation.

596 **Penalty for blocking**

$$c = \begin{cases} 0.2, & \text{successfully complete this action} \\ 0, & \text{cannot complete this action in 0.5s} \end{cases} \quad (1)$$

597 **Penalty for Slippery**

$$s = \begin{cases} 0.5, & ||y_t^i - y_t^{i-1}|| \geq 0.5 \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

598 The y_t^i and y_t^{i-1} stands for the embeddings of the tactile images in the current step i and the last step
599 $i - 1$.

600 **Overall Rewards**

$$R = \begin{cases} 1, & \text{if success} \\ \alpha D_{KL}(P||Q) + \beta d + \gamma \cdot o - c - s, & \text{otherwise} \end{cases} \quad (3)$$

601 In Eqn. 3, P is the executed trajectory, Q is the expert trajectory, $D_{KL}(P||Q)$ is the KL Divergence
602 between the executed trajectory and the expert trajectory, d, o, c, s are defined above. The setup of
603 each weight: $\alpha = 0.5, \beta = 0.3, \gamma = 0.2$.

604 E Emulate the Physical Environment for Policy Evaluation

605 To emulate the physical robot environment, we introduce random noise to those 10 unseen data
606 sequences. The robot state space input undergoes a random position noise within the range
607 $[-3\text{cm}, +3\text{cm}]$ for each axis. Gaussian noise, denoted as $\mathcal{N}(0, \sigma)$, is added to both the tactile image
608 and audio signal. In this notation, $\mathcal{N}(0, \sigma)$ signifies a Gaussian distribution with a mean of 0 and a
609 standard deviation of σ . For tactile images, the noise affects pixel values in the range $[0, 255]$, while
610 for audio data, it impacts signal values in the range $[0, 1]$. Given their distinct ranges, we apply
611 Gaussian noise with standard deviations of $\sigma = 100$ for tactile images and $\sigma = 0.4$ for audio data.

612 F MSE Loss

613 For calculating the MSE Loss between two action sequences, we need to normalize the actions'
614 translation vectors and rotation vectors since they have different scales. Specifically, we use min-
615 max normalization on both the translation vectors and rotation vectors, where the max vector and
616 min vector are selected from the training set. As a result, translation vectors and rotation vectors
617 will have the same scale for calculating the MSE Loss. The formula is shown as:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

618 Where: y_i represents the ground truth normalized action, \hat{y}_i represents the generated normalized
619 action, and n is the number of all action steps.

620 G Ablation Study: Do Multi-Modal Tactile Feedback Improve the Task 621 Performance?

622 In this section, we evaluate the performance of our multi-modal tactile embeddings. We consider
623 the following baselines: i). *MimicTouch w/o T & A*: MimicTouch without tactile or audio embed-
624 dings, ii). *MimicTouch (T)*: MimicTouch incorporating only tactile embeddings, iii). *MimicTouch*

Models	<i>MimicTouch w/o T & A</i>	<i>MimicTouch (T)</i>	<i>MimicTouch (A)</i>	<i>MimicTouch (T + A)</i>
MSE Loss	0.62	0.38	0.48	0.21
Success Rate	4%	24%	16%	40%

Table 3: MSE Loss over test sets and Task success rates of 25 policy rollouts.

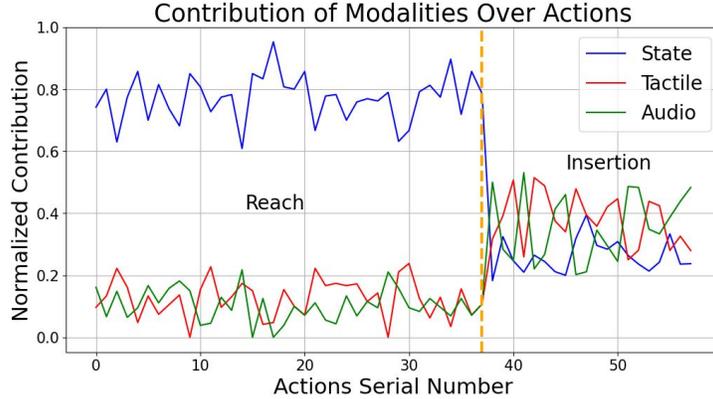


Figure 7: Visualization of the impact of each sensor modality during policy execution. The left side of the dashed orange line is Reach Phase, and the right side is Insertion Phase.

625 (A): *MimicTouch* incorporating only audio embeddings., and iv). *MimicTouch (T + A, Ours)*: *MimicTouch*
626 incorporating both tactile and audio embeddings. We evaluate the policy performance in
627 terms of the MSE losses over the test sets (see Sec. 4.2.1) and task success rates over 25 policy
628 rollouts. In addition, we visualize the impact of each sensor modality during policy execution.
629 Specifically, we plot the normalized distance between the query feature and the selected key feature
630 for each sensor input. A larger distance means that the corresponding sensor modality contributes
631 more in selecting the key feature from the demonstration library.

632 From Table. 3, we observe that without using both tactile images and audio signals, the MSE loss
633 (task success rate) is 0.62 (4%), which is significantly higher (lower) than the others. By incorporat-
634 ing both tactile and audio feedback, the MSE loss can be as low as 0.21, and more importantly, the
635 task success rate can reach 40%.

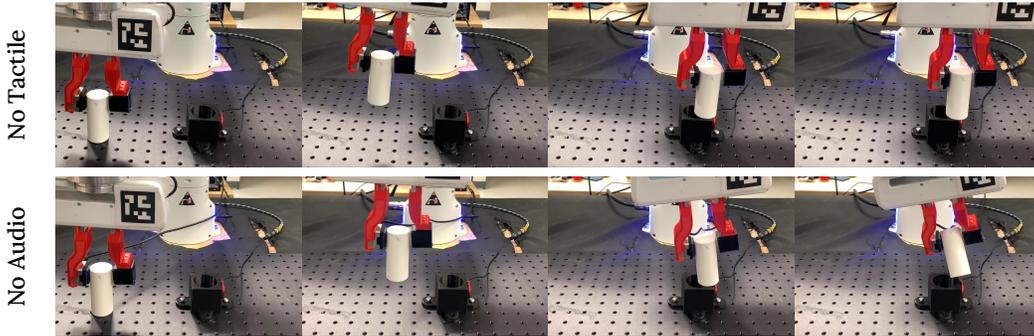


Figure 8: Policy rollouts of some failure examples with only tactile feedback or only audio feedback.

636 As shown in Fig. 7, tactile and audio inputs start to play important components during the Insertion
637 phase, when most contacts occur. We also have qualitative results shown in the Fig. 8. According
638 to those results, we can find that a lack of tactile feedback easily leads to incorrect motion when
639 contact appears, whereas the lack of audio feedback easily leads to an inability to detect external
640 collisions.

641 Therefore, we can conclude that the multi-modal tactile feedback is crucial for the success of
642 contact-rich insertion tasks.

643 H Generalization Setting

644 In this section, we describe the setting of each generalization task.

645 **Shifting Positions:** In this generalization task, we shift the hole positions for 0.8 cm in either $\pm x$
646 or $\pm y$ to test the generalization ability for finishing the insertion task with under varied alignment
647 conditions.

648 **Tilting Angles:** In this generalization task, we tilted the hole angle for 10 degrees or 20 degrees to
649 test the generalization ability for finishing the insertion task with different contact positions.

650 **Two-stage Dense Packing:** We introduce the two-stage dense packing task, which requires the robot
651 to perform two consecutive alignment adjustments to complete the dense packing process. Each hole
652 will challenge the robot’s ability to adjust its alignment according to tactile feedback efficiently.

653 **Multi-material Dense Packing:** In this generalization task, the robot is required to insert the cylinder
654 into the hole which contains multiple objects: pen, tissue, and sponge. This setting has rigid
655 objects (pen) and deformable soft objects (tissue and sponge) with different materials and shapes,
656 which challenge the robot’s ability to accomplish the task with different tactile feedback from dif-
657 ferent materials.

658 **Furniture Assembly:** In this generalization task, the robot is required to insert the cylinder into a
659 small hole in a table for screwing. This task will test two aspects of our policy: whether the robot
660 can insert the object into a smaller hole, and whether the robot can adjust it to a position, that is deep
661 enough to be skewed by a human-defined simple script (to rotate the end-effector for 120°), based
662 on the tactile feedback from the threads in the hole.

663 **Policy Evaluation Process:** We evaluate the policy performance in those five different generaliza-
664 tion settings for both MimicTouch final policy and the *Openloop* Policy. For *Shifting Positions*, we
665 rollout the policies 10 times in each direction of $\pm x$ or $\pm y$, resulting in a total of 40 evaluations. For
666 *Tilting Angles*, we rollout the policies 25 times for both tilting directions in $\pm x$, resulting in a total
667 of 50 evaluations for the 10° tilting and 20° tilting, respectively. For *Two-stage Dense Packing*, we
668 rollout the policies 25 times. For *Multi-material Dense Packing*, we rollout the policies 25 times on
669 both rigid object “pen” (Rigid in the table) and deformable soft objects “tissue and sponge” (Soft in
670 the table). For *Furniture Assembly* (Assem in the table), we rollout the policies 25 times. This task
671 has two sub-evaluation metrics: insertion (Assem (I) in the table) and adjustment (Assem (A) in the
672 table). Notably, the success rate for Assem (I) is the success rate for the number of attempts (25),
673 and the success rate for Assem (A) is the success rate when the insertion is successful.

674 I Generalization Results

675 In this section, we summarize the zero-shot generalization results based on the quantitative results
676 shown in Table. 2, and the qualitative results shown in Fig. 9.

- 677 • MimicTouch policy can zero-shot transferring to insertion tasks with different contact po-
678 sitions, tilted angles, and even different sizes of holes (Asembly (Insertion)).
- 679 • For the Two-stage Dense Packing, MimicTouch policy displays its robustness to adjust
680 according to multiple stages of contact information according to the quantitative result and
681 the video. This shows that our model can make continuous and correct adjustments based
682 on the continuously varied contact information.
- 683 • MimicTouch policy also shows its power in the multi-material task domains. Due to dif-
684 ferent materials in the environment, the sensor feedback will be different from the training
685 environment. In this case, our policy still has great performance on other rigid objects (pen).
686 For the deformable soft object (tissue and sponge), The success rate is a bit lower because
687 of two challenging issues: i). it’s hard to get audio feedback for contact with soft objects,
688 ii). sponge sometimes is too soft to get tactile feedback. With those issues, our policy still
689 gets 64% success rate. Moreover, the qualitative result from the video shows impressive
690 performance in adjusting continuously according to the deformation of the tissue.

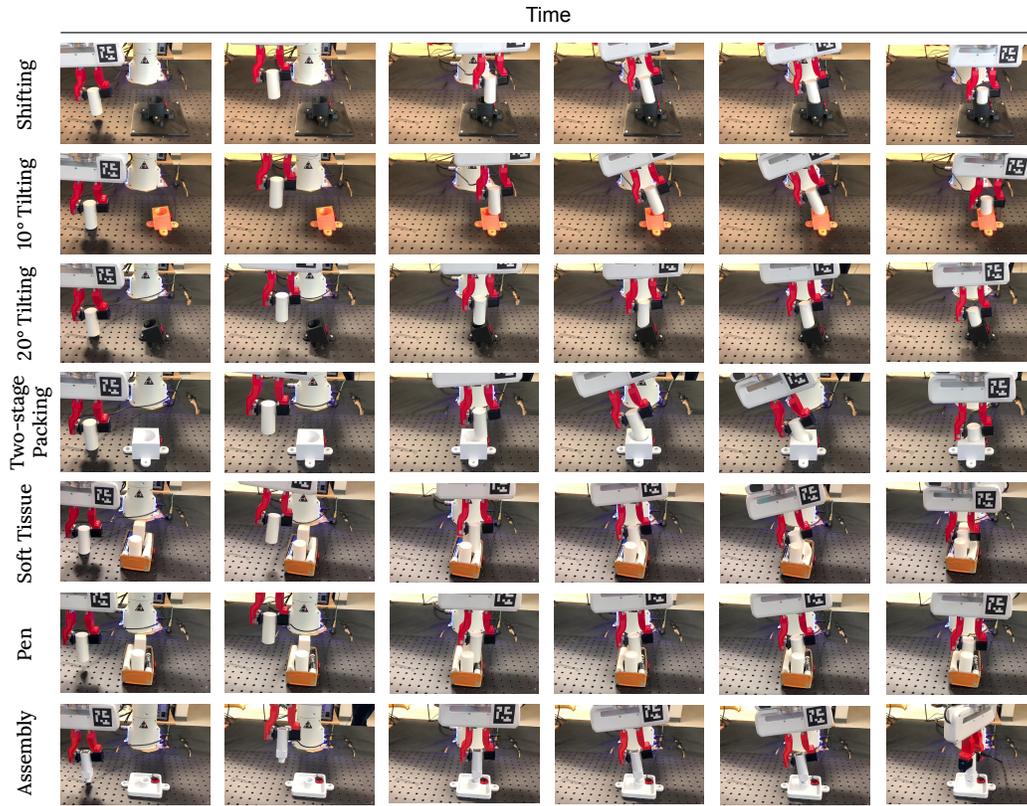


Figure 9: Task setup and qualitative results for zero-shot generalization tasks

691
692
693
694

- In the assembly task, MimicTouch policy not only can insert the object into a small hole but also can adjust the object to a correct pose and insert it to a deep-enough position according to the tactile feedback from the contact between screw threads. This allows us to use a very simple script (to rotate the end-effector for 120°) to solve the assembly task.