

# Data and Learning Where it Matters for Contact-Rich Manipulation

Anonymous Author(s)

Affiliation

Address

email

1       **Abstract:** Learned policies trained end-to-end on large datasets often remain brittle  
2       in high-precision tasks and struggle with generalization. We find that these  
3       limitations largely stem from a lack of structure and focus in data collection. Our  
4       key insight is to leverage dense data collection only for the critical segment of  
5       contact-rich tasks and to rely on traditional planning during simple free-space  
6       motion. We propose an automated data-collection scheme in combination with  
7       offline deep reinforcement learning for the critical segment of the task, eliminating  
8       reliance on a teleoperator’s skill and on online policy updates. Across four  
9       challenging real-world tasks, using only 2–2.5 h of autonomous data collection,  
10      we achieve an average success rate of 96 %, compared to the strongest baseline at  
11      55 %. Notably, performance remains high in out-of-distribution scenarios where  
12      end-to-end approaches struggle. Our results pave the way for targeted data collection  
13      for contact-rich tasks and for high success rates in precision applications. *We*  
14      *will upload all our videos, training datasets, and evaluation datasets. Website.*

15      **Keywords:** Data for Robotics, Learning, Manipulation, Foundation Models, DRL

## 16   1 Introduction

17   End-to-end robot learning has recently become very popular [1, 2, 3], especially in contact-rich  
18   manipulation tasks that are hard to model analytically. This progress has been further accelerated  
19   by large-scale data collection [4, 5]. Yet certain limitations of end-to-end policies become apparent:  
20   they struggle with high-precision tasks, and generalization to out-of-distribution (OOD) scenarios  
21   remains limited. It is still unclear whether simple data scaling will close this gap, especially given  
22   that data collection pipelines remain expensive [6, 7].

23   We argue that at least part of this gap is not merely a matter of scale, but of *what* data is collected  
24   and *how* tasks are structured. Specifically, we show that high-precision tasks usually fail at the  
25   contact-rich, critical segment — such as an insertion with tight clearance. The other parts of the  
26   task typically involve unconstrained free-space motion and are thus easier to handle. End-to-end  
27   approaches treat these fundamentally different phases uniformly during data collection and policy  
28   learning [1, 2, 3, 6, 7].

29   We propose to exploit this structure explicitly. Our key insight is to *densely collect data where it*  
30   *matters* — at the critical segment. For the remainder of the task, we go to the other extreme and  
31   solve it without using robot-specific data, using off-the-shelf pose estimation and motion planning.  
32   Recent advances in language-guided image segmentation [8] and pose estimation [9] make such  
33   sequential pipelines increasingly practical for real-world manipulation while reducing the reliance  
34   on robot data.

35   First, we introduce a scheme for autonomous data acquisition that performs dense data collection  
36   at the critical segment, in combination with offline deep reinforcement learning (DRL) for policy  
37   learning. The autonomous data-acquisition scheme eliminates the constant dependency on a human  
38   teleoperator. For deployment, we propose combining the learned policy with motion planning for

Submitted to the 10th Conference on Robot Learning (CoRL 2026). Do not distribute.

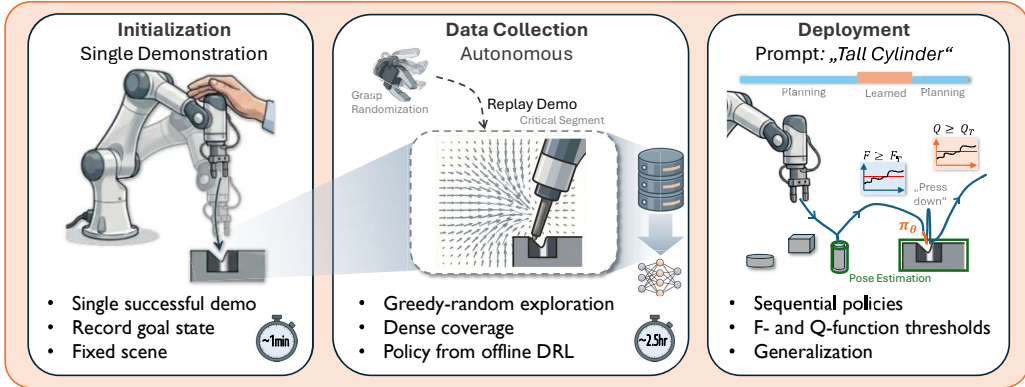


Figure 1: **Method.** (Left) First, we record a single demonstration in the scene using kinesthetic teaching. The scene layout remains unchanged for the subsequent data collection. (Middle) The demonstration is replayed until the critical segment is reached, and then we execute a mixed greedy-random policy  $\pi_{\text{explore}}$  for dense data collection, yielding successful and failed rollouts. We combine data collection with offline DRL to obtain a learned policy  $\pi_{\theta}$ . (Right) In deployment,  $\pi_{\theta}$  for the contact-rich segment is sequentially chained with off-the-shelf pose estimation and motion planning.

39 free-space motion and switching policies based on contact events and Q-function values. We evaluate  
 40 our approach on four challenging real-world tasks and difficult out-of-distribution scenarios.  
 41 Using only 2 h to 2.5 h of wall-clock time for autonomous data collection, our method achieves  
 42 success rates of 94 to 98 % across four real-world tasks and 96 % on average, outperforming base-  
 43 lines that achieve no more than 55 %. While end-to-end methods struggle under out-of-distribution  
 44 conditions, our approach maintains high success rates.

45 In summary, we present three core contributions. (1) We propose a compositional framework that  
 46 leverages dense data collection and a learned policy in an offline manner for the critical task segment,  
 47 while relying on motion planning for free-space motion. (2) We introduce an autonomous data-  
 48 collection scheme that leverages mixed greedy-random exploration for dense data collection at the  
 49 critical segment. (3) During deployment, we integrate planning and learning using contact events  
 50 and Q-functions for policy switching, achieving success rates above 94% across all tasks.

## 51 2 Related Work

52 **Data collection strategies.** Imitation Learning typically relies on teleoperation data [2, 10, 11].  
 53 DRL on hardware similarly requires human intervention [12]. Automated data acquisition schemes  
 54 are mostly limited to simulation. Assembly-by-disassembly generates demonstrations by automatic  
 55 disassembly and trains a DRL policy using imitation rewards [13]. Other approaches use search for  
 56 data generation [14, 15] or to augment initial human demonstrations [16]. In contrast, we propose  
 57 an automated data-acquisition scheme for the real world for the critical task segment.

58 **Contact-rich manipulation.** Methods relying on pure motion planning use fine-tuned pose esti-  
 59 mation to achieve high success rates [17, 18, 19]. Learned policies achieve close to 100% when  
 60 ground-truth 6D object poses are available [13, 20, 21]. HIL-SERL demonstrated 100% success rate  
 61 with DRL and vision input, but only for short-horizon tasks [12, 22]. Imitation policies typically  
 62 achieve around 80% success rate on complex tasks [2, 23, 24], and we show that their failure cases  
 63 are concentrated at the critical segment. For traditional control methods, direct force or hybrid force-  
 64 position control has been proven essential for contact-rich tasks [25, 26], with first learning-based  
 65 models adopting the approach [3, 27, 28, 29].

66 **Modular Policies.** A large body of work separates different stages of manipulation tasks into ded-  
 67 icated sub-policies [30, 31], which can also be trained independently on heterogeneous data and  
 68 combined through diffusion guidance [32] or a learned router [33]. For robotic foundation models,  
 69 mixture-of-experts architectures have become popular, enabling scalable multitask learning by spe-

70 cializing lightweight experts across behaviors [34, 35]. In contrast to prior approaches that compose  
 71 multiple learned policies, we exploit the structure of contact-rich manipulation itself, learning only  
 72 the critical segment to improve robustness and OOD generalization.

### 73 3 Method

#### 74 3.1 Problem Setup

75 We model contact-rich interaction as an MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \rho, \mathcal{P}, r, \gamma)$ , where  $\mathbf{s}_t \in \mathcal{S}$  and  $\mathbf{a}_t \in \mathcal{A}$   
 76 are the state and action at time step  $t$ ,  $\rho(\mathbf{s}_0)$  is a distribution over initial states,  $\mathcal{P}(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$   
 77 are the unknown transition dynamics,  $r_t = r(\mathbf{s}_t, \mathbf{a}_t) \in \mathbb{R}$  is a reward function encoding the task,  
 78 and  $\gamma \in [0, 1)$  is a discount factor. We observe the scene with a single wrist camera and obtain  
 79 state estimates of task-relevant objects using pose estimation. Our goal is to collect a dataset for  
 80 training a policy  $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$  for the critical segment parameterized by a neural network with  
 81 parameters  $\theta$  that maximizes the expected discounted return  $J(\pi_\theta) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t]$ . An overview  
 82 of our methodology is provided in Figure 1.

#### 83 3.2 Autonomous Offline Data Collection

84 We first set up the scene with the objects at fixed locations in the workspace. If necessary, fixtures can  
 85 be 3D printed to keep the objects in place (Appendix A.1). Then, we record a single demonstration  
 86 that performs the task using, for example, kinesthetic teaching, and put all objects back to their initial  
 87 locations. Our further scheme relies on keeping the scene layout unchanged during data collection.

88 We start by replaying the demonstration trajectory until the critical segment of the task is reached.  
 89 In practice, we define the critical task segment based on the uncertainty of the pose estimation  
 90 (Appendix A.5). Then, we execute an explorative data-collection policy

$$\pi_{\text{explore}}(\mathbf{s}_t) = \mathbf{a}_t = \begin{cases} \mathbf{a} \sim \mathcal{U}(\mathcal{A}) & \text{with probability } p \\ \mathbf{a}_t^* & \text{with probability } 1 - p, \end{cases} \quad (1)$$

91 where  $\mathcal{U}$  is the uniform distribution,  $p \in [0, 1]$  is a probability threshold, and  $\mathbf{a}_t^*$  is a greedy action.  
 92 There are multiple ways to instantiate greedy policies [36, 37, 38, 13]. We use the action that  
 93 minimizes the distance to the goal state  $\mathbf{s}_g$  with a velocity  $v$ :

$$\mathbf{a}_t^* = v \cdot \frac{(\mathbf{s}_g + \Delta \mathbf{s}_{\text{grasp}}) - \mathbf{s}_t}{\|(\mathbf{s}_g + \Delta \mathbf{s}_{\text{grasp}}) - \mathbf{s}_t\|_2}, \quad (2)$$

94 where  $\Delta \mathbf{s}_{\text{grasp}}$  is the randomized grasp offset to the reference demonstration, and the state is the  
 95 Cartesian end-effector pose. Additionally, we apply a safety filter to the action (1) to prevent explo-  
 96 ration from drifting outside the critical segments, constraining the end-effector position to a sphere  
 97 with center  $\mathbf{s}_{\text{safe}}$  and radius  $r_{\text{safe}}$ .

98 This combination of random and optimal actions provides dense coverage of the critical segment  
 99 while maintaining high success rates. The data collection episodes end either by truncation (timeout-  
 100 based) or by task completion. The task reward  $r_t$  is a sparse reward determined by the distance of  
 101 state  $\mathbf{s}_t$  to the final state of the hardcoded trajectory  $\mathbf{s}_g$ :  $r_t = R_{\text{success}} \cdot \mathbb{1}(\|\mathbf{s}_t - (\mathbf{s}_g + \Delta \mathbf{s}_{\text{grasp}})\|_2 \leq$   
 102  $\epsilon)$ . Computing the reward state-by-state eliminates the need for a pre-trained vision-based reward  
 103 classifier [39, 40], which would require additional demonstrations. Instead, we obtain a vision-based  
 104 classifier as the result of the DRL training in the form of the Q-function.

105 Contact-rich manipulation requires tight integration of perception and force feedback, particularly  
 106 during the insertion phase [29]. For this reason, we employ hybrid force-position control [25, 29]  
 107 during the critical, force-sensitive interactions. Specifically, we apply a desired force along the  
 108 insertion axis and position-control the remaining degrees of freedom. This enables compliant and  
 109 force-controlled behavior with precise positioning, and the same is applied during deployment.

110 **3.3 Policy Learning**

111 To minimize data collection effort, we require a sample-efficient learning algorithm. Similar to  
 112 RLPD [41], we build on top of SAC [42] and employ improvements for sample efficiency. To mit-  
 113 igate Q-function overestimation, we add layer normalization [43] to the critic and use Randomized  
 114 Ensembled Double Q-Learning [44], where we train  $N_Q$  Q-functions and select a random subset  $S$   
 115 of size  $N_{Q,S}$  to compute the TD-targets as

$$y_t = r_t + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_\theta} \left[ \min_{i \in S} Q_{\phi_i}^-(\mathbf{s}_{t+1}, \mathbf{a}') - \alpha \log \pi_\theta(\mathbf{a}' | \mathbf{s}_{t+1}) \right]. \quad (3)$$

116 The Q-functions and policy are trained via the standard SAC objective [42] using  $y_t$ . Since our  
 117 setting does not rely on online learning, we fix the entropy coefficient  $\alpha$  rather than adapting it  
 118 during learning. This setup is a significant algorithmic and architectural simplification compared to  
 119 online policy optimization.

120 **3.4 Full Task Execution**

121 Interaction with the scene is done using a language-guided segmentation model (SAM3 [8]), which  
 122 we prompt with commands such as "purple Lego brick", "gray plastic cover with circular grille",  
 123 and "yellow cardboard salt box" for our tasks. We feed the segmentation masks into an off-the-shelf  
 124 pose estimation model [9, 45] to obtain 6D pose estimates for the two objects in the scene that we  
 125 manipulate, unless otherwise specified. Pose estimation is run multiple times to refine the estimates  
 126 throughout the trajectory (Appendix A.5). The pose estimates serve as the basis for motion planning,  
 127 which uses simple waypoint following. Our formulation can easily be extended to include trajectory  
 128 optimization or obstacle avoidance.

129 All waypoints for motion planning are defined relative to the object’s 6D pose and based on the single  
 130 demonstration. The robot first moves towards the pick-up object until it reaches a specified waypoint.  
 131 Then, we grasp the object and transport it to the critical segment (also using waypoints). The learned  
 132 policy is triggered once contact is made, i.e., if  $\|\mathbf{F}\|_2 \geq \mathcal{F}_{\text{threshold}}$  [46, 47] (Appendix A.2). Then, the  
 133 learned policy is executed and terminated based on the learned Q-functions and a threshold  $\lambda_{\text{success}}$   
 134 to detect successful insertion:

$$Q(\mathbf{s}_t, \pi_\theta(\mathbf{s}_t)) = \frac{1}{N_Q} \sum_{i=1}^{N_Q} Q_{\phi_i}(\mathbf{s}_t, \pi_\theta(\mathbf{s}_t)) > \lambda_{\text{success}}, \quad (4)$$

135 We additionally use a hysteresis-style detection to avoid false negative cases. Further details and ex-  
 136 amples for the Q-function-based success classification are in Appendix A.4. After the learned policy  
 137 terminates, any remaining motions, such as a push-down to complete the insertion, are executed.

138 **4 Experimental Results**

139 **4.1 Task descriptions**

140 We evaluate our approach on multiple challenging real-world tasks that require handling deformable  
 141 components, precise alignment, and contact-rich interaction, as shown in Figure 2. Experimental  
 142 details are provided in Appendix A.7.

143 **Shelf stocking:** This task requires stocking a tightly packed shelf with a yellow salt box. Most  
 144 objects on the shelf are made of cardboard and may deform during interaction. The policies need to  
 145 learn to react to those changes, precisely align the salt box, and perform an additional push to fully  
 146 complete the insertion. *Partial success* is if the object is placed on the shelf, but not inserted.

147 **Lego stacking** tasks have recently enjoyed popularity in the robot learning community. The task is  
 148 to pick up a brick and stack it on another  $4 \times 2$  brick from an unknown position on the base plate.

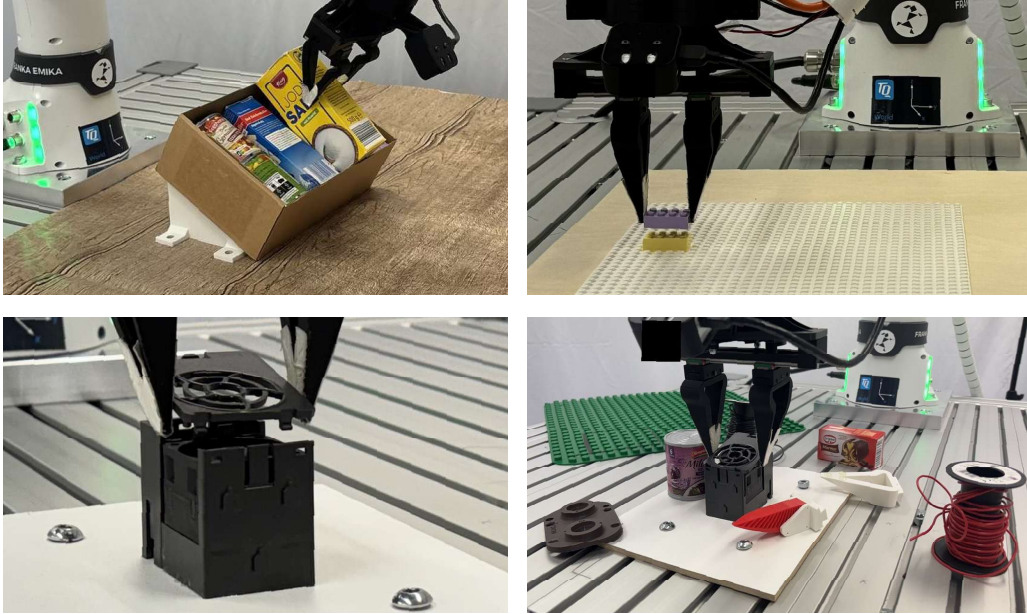


Figure 2: **Tasks.** We evaluate our method on four challenging real world tasks: shelf stocking, Lego stacking, mounting two PBT/PC-molded parts (fan cover), and a difficult version of the fan cover task that we train on scene distractors, similar objects, and changing ground surfaces. For all tasks, we evaluate the policies on the additional out-of-distribution scenarios shown in Appendix A.6.

149 The task is solved once the two bricks are fully stacked together. *Partial success* is defined as the  
 150 blocks being placed on top of each other, but not fully inserted.

151 **Fan cover** originates from a real-world production line that currently necessitates manual assembly.  
 152 The assembly requires aligning the studs through intricate, contact-rich motions against the base.  
 153 The injection-molded PBT/PC are deformable in the process. *Partial success* is defined as the cover  
 154 being placed on the base, but not fully inserted. **Fan cover (difficult):** A difficult version of this task  
 155 introduces additional diversity during data collection and deployment (Figure 2). For both versions  
 156 of this task, we assume the fan base is fixed and do not estimate it.

## 157 4.2 Implementation Details

158 For our data collection scheme, we set the probability of random actions to  $p = 0.8$  and truncate  
 159 episodes after 10 s (150 steps). We train our policies for 5 epochs, which takes 1 h on a GPU  
 160 workstation. The policy receives a single RGB viewpoint from a wrist-mounted camera, and pose  
 161 estimation uses additional depth measurements [9, 45]. The images are cropped to the critical task  
 162 segment (Appendix A.9), embedded using DINOv2 [48], and projected down to 16 dimensions using  
 163 an MLP before concatenation with the state. The state includes F/T measurements, controller error,  
 164 Cartesian velocity, and the previous action. During deployment, we use  $\lambda_{\text{success}} = 0.93R_{\text{success}}$  for  
 165 the Q-function based success classifiers. Further hardware and training details are in Appendix A.8.

166 **Baselines.** We benchmark against end-to-end Imitation Learning and DRL. For Imitation, we choose  
 167 DiTFlow [49] and DP [1] as task-specific policies and train them from scratch for 50k gradient steps  
 168 (5 h). For DiTFlow, we include a version trained on data collected by a novice teleoperator inexpe-  
 169 rienced with data collection.  $\pi_{0.5}$  [2] is our baseline for a foundation model that we finetune for 40k  
 170 steps on  $8 \times \text{H100 GPUs}$  (18 h). For DRL, we use HIL-SERL [12], with 20 initial demonstrations  
 171 and a human guiding policy learning as per the original implementation. We include a baseline with  
 172 pure pose estimation and motion planning (PE & MP) without any learning.

Table 1: Data collection

Task	# Episodes (Failures <sup>1</sup> ) / # Datapoints / Data collection wall-clock time			
	Teleop	Teleop (Novice) <sup>2</sup>	HIL-SERL	Ours
Shelf stocking	100 (22) / 24k / 2.3 h	100 (18) / 34k / 2.3 h	346 (133) / 27k / 3.2 h	<b>300 / 27k / 2 h</b>
Lego stacking	100 (30) / 33k / 2.5 h	100 (27) / 42k / 2.6 h	495 (221) / 32k / 3.2 h	<b>500 / 30k / 2.5 h</b>
Fan Cover	100 (50) / 39k / 2.75 h	—	—	<b>300 / 33k / 2 h</b>
Fan Cover (hard)	100 (55) / 40k / 2.5 h	—	—	— <sup>3</sup>

<sup>1</sup>Episodes that failed during teleoperation — i.e., where the teleoperator did not successfully complete the task — are indicated in brackets and were excluded from policy training. <sup>2</sup>“Novice” refers to an inexperienced teleoperator who has been given a couple of trials before starting data collection. All other data have been collected by an experienced operator. <sup>3</sup>No additional data collected.

Table 2: Success Rates (Partial Success Rates<sup>1</sup>) [%]

Task	DiTFlow	DiTFlow (Novice) <sup>2</sup>	DP	Finetuned $\pi_{0.5}$	MP & PE <sup>3</sup>	HIL-SERL	Ours
Shelf stocking	38 (96)	34 (80)	44 (100)	42 (92)	98 (98)	10 (24)	<b>98 (100)</b>
Lego stacking	60 (88)	22 (48)	44 (92)	6 (80)	78 (100)	6 (22)	<b>94 (100)</b>
Fan Cover	30 (88)	—	30 (96)	20 (96)	22 (94)	—	<b>96 (98)</b>
Fan Cover (hard)	18 (82)	—	26 (92)	12 (76)	22 (94)	—	<b>94 (96)<sup>4</sup></b>
Average	37 (89)	—	36 (95)	20 (86)	55 (97)	—	<b>96 (99)</b>

We roll out each policy for 50 trials. <sup>1</sup>Partial success is when the objects are placed correctly but not fully inserted or assembled. <sup>2</sup>“Novice” refers to a policy trained on data collected by a novice teleoperator. All other data have been collected by an experienced operator. <sup>3</sup>Motion Planning & Pose Estimation. <sup>4</sup>We deployed the same policy as for the simple version of the task without additional data collection.

### 173 4.3 Policy Evaluations

174 Our central claim is that traditional planning combined with data collection and learning for the  
 175 critical segment leads to high success rates and robust policies. We aim for equal wall-clock time  
 176 for data collection across all methods (Table 1). During the same wall-clock time, our data collection  
 177 scheme collects significantly more data at the critical segment. We note that despite being largely  
 178 autonomous, our data collection scheme required a few interventions when the scene was not reset  
 179 correctly. Those interventions put significantly less strain on the operator than teleoperation.

180 While end-to-end methods capture the task semantics, they typically struggle at the critical seg-  
 181 ment, as shown by high *partial* success rates — indicating correct grasping and coarse placement  
 182 — but low overall success rates (Table 2). HIL-SERL struggled primarily with undirected gripper-  
 183 actions during exploration, leading to frequent failures given the long task horizon. We also observe  
 184 high variance in training outcomes during runs, which is characteristic of online policy optimiza-  
 185 tion [50, 51]. Our sequential pipeline, consisting of MP, PE, learned policy, and policy switching,  
 186 achieves consistently  $\geq 90\%$  SR. Remaining failure cases are due to OOD scenarios and hardware  
 187 inaccuracies. See Appendix A.10 for an overview of failure cases.

188 A modular approach enables tailoring policy parameters for each task segment. For instance, during  
 189 the planned push-down motion, we control the force to a magnitude similar to that in the single  
 190 kinesthetic demonstration. This leads to overall fewer forces being applied to the objects, which is

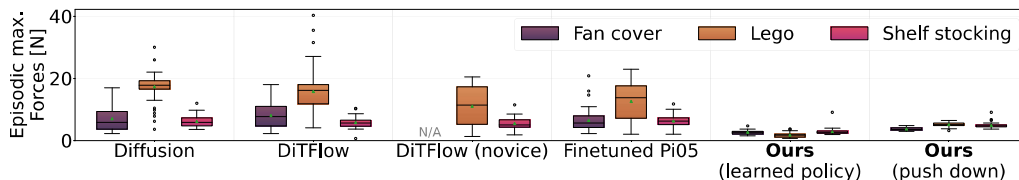


Figure 3: Maximum forces applied during policy rollouts. **Ours** applies significantly less forces and torques to the objects than baselines. In the figure, we divide ours between the learned insertion policy and the planned push-down motion. Torques are shown in Appendix A.12.

Table 3: Success Rates (Partial Success Rates<sup>1</sup>) [%] for out-of-distribution scenarios.

Task	DiTFlow	DP	Finetuned $\pi_{0.5}$	HIL- SERL	Ours
Shelf stocking	20 (40)	25 (45)	10 (65)	0 (5)	<b>90 (100)</b>
Lego stacking	10 (10)	10 (25)	0 (0)	0 (0)	<b>85 (95)</b>
Fan Cover	0 (30)	0 (20)	0 (50)	—	<b>90 (95)</b>
Fan Cover (hard)	0 (50)	5 (50)	5 (50)	—	<b>95 (95)<sup>2</sup></b>

We roll out each policy for 20 trials. <sup>1</sup>Partial success is when the objects are placed correctly but not fully inserted or assembled. <sup>2</sup>We deployed the same policy as for the simple version of the task — without additional data collection for this task.

191 critical for the sensitive cardboard and PBT/PC molded parts (Figure 3). Applied maximum forces  
192 are reduced on average by 49% ( $-4.7$  N) compared to baselines.

193 A policy trained on data from a novice operator (DiTFlow (Novice)) achieves lower success rates  
194 (Table 2), demonstrating the dependency on demonstration data and motivating automated data  
195 collection. We find that novice demonstrations are characterized by multi-modal demonstrations,  
196 i.e., the task is demonstrated in different ways, higher speed variance ( $CV_{\text{novice}} = 0.285$  vs.  
197  $CV_{\text{expert}} = 0.185$ ), and overall longer trajectories (27 s vs. 24 s).

#### 198 4.4 Generalization

199 We expose the policies to OOD scenarios, including distractor objects, different object configu-  
200 rations, placements, and backgrounds (Appendix A.6). End-to-end methods struggle with OOD  
201 scenarios (Table 3). The fan cover (hard) task includes the standard fan cover OOD scenarios dur-  
202 ing training, and during evaluation additional distractors are introduced. The augmentations during  
203 training increased partial success rates to 50%, but fully complete insertions remain few (5%). Sub-  
204 stantially more diverse data would likely be required to achieve generalization.

205 As our policies operate on local camera images, the success rates stay similar to in-distribution eval-  
206 uations. Remaining failure cases include distractor objects being placed close to the insertion point,  
207 taking local visual observations out of distribution (Appendix A.10). We present additional qualita-  
208 tive results for complex scenarios such as dynamically grasping from a human hand, arbitrary object  
209 poses, and further randomizations — all without additional data collection — in Appendix A.6. The  
210 motion planning can easily be extended to include trajectory optimization and avoidance [52].

#### 211 4.5 Data Collection Ablations

212 We analyze key parameters of our autonomous data acquisition scheme. We use a MuJoCo [53]  
213 simulated setup of the Lego tasks, train policies on three different seeds, and report the results on  
214 1000 policy rollouts. The results focus on three insights.

215 **Successful policies require dense data coverage of the critical segment (Figure 4 a+b).** Policies  
216 trained on data with little exploration ( $p \leq 0.3$ ) achieve low success rates. Little exploration  
217 results in short trajectories that yield little training data (see Figure 4 (a)). Note that  $p = 0$  is close  
218 to imitation learning data. On the other hand, pure exploration ( $p = 1$ ) leads to low SR and thus  
219 unsuccessful policies. In practice, one requires sufficiently long rollouts while maintaining a high  
220 success rate, which we observe for  $0.6 \leq p \leq 0.85$ . The time spent at the critical segment per  
221 rollout is critical, necessitating DRL to learn from non-optimal demos.

222 **Scaling of success rate and dataset size with pose estimation uncertainty (Figure 4 c).** Increas-  
223 ingly uncertain pose estimation expands the critical segment of the task. In simulation, we require  
224 400 rollouts for 6 mm and 700 rollouts for 12 mm uncertainty. This corresponds to 2 h to 3.5 h of  
225 data collection in the real world. We tend to require more data in simulation than on hardware for  
226 the same success rates, possibly because of inconsistent simulation physics during contacts.

227 **Multi-modal sensing for high success rates (Appendix A.11).** Combining a wrist-mounted F/T  
228 sensor with a single wrist camera achieves 100% success in simulation. This is the setup we also

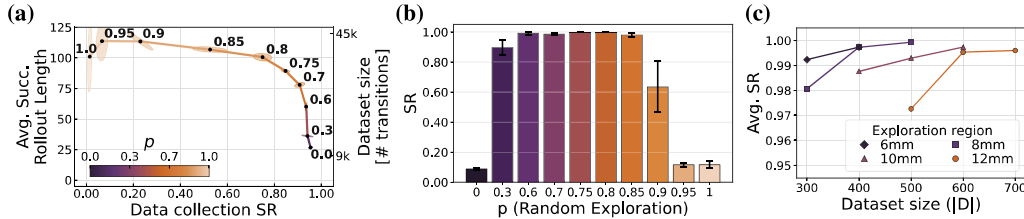


Figure 4: Ablation of data collection parameters. **(a)** Effect of random action sampling rate  $p$  on average successful rollout length, dataset size, and success rate during data collection. **(b)** Trained policy success rate for different random exploration rates  $p$  during data collection. **(c)** Scaling of SR and dataset size with pose estimation uncertainty.

229 chose for our hardware experiments. F/T measurements provide useful feedback even when the task  
 230 state is not fully observable from images alone, e.g., when objects are already in contact.

## 231 5 Discussion

232 We use motion planning for the free-space motion of the task, but this doesn't necessarily have to be  
 233 the case. Policies for free-space motion could come from other inexpensive, non-robot data sources  
 234 that can be collected at a large scale, such as human videos [54]. Future work can investigate how  
 235 to balance, collect, and integrate different data sources for modular policies.

236 Sequential policies naturally give rise to task-segment specific priors, which increase performance  
 237 significantly for difficult tasks, as we demonstrate with reduced contact forces compared to end-  
 238 to-end methods. A key question is the design of switching mechanisms, and Q-function-based  
 239 switching worked well for our tasks. VLM-based task segment switching for semantically-rich  
 240 tasks is an interesting direction for future research.

241 All vision-based policies are sensitive to visual OOD shifts. However, locally operating policies  
 242 narrow down the set of OOD scenarios and remain robust under scene-level distractors. In contrast  
 243 to end-to-end methods, visual randomization can remain tractable for local policies by leveraging  
 244 simulation and sim-and-real cotraining [55] — potentially without additional real robot data.

## 245 6 Conclusion

246 Robot data is expensive, and we should consider how to collect and use it efficiently. While clearly  
 247 more robot data is needed overall, we argue for focusing collection on tasks and task-segments  
 248 that are otherwise difficult to solve. By adopting a modular approach with targeted data collection,  
 249 we maintain high success rates for OOD scenarios, demonstrating that targeted reliance on robot  
 250 data and structured task decomposition are key enablers of robust generalizable manipulation while  
 251 keeping data acquisition efforts at bay.

## 252 7 Limitations

253 **Task diversity.** We demonstrate high success rates on tasks with sequential separation of free-space  
 254 and contact-rich motion. For other tasks, such as t-shirt folding, sequential boundaries are harder to  
 255 define — motivating future work in policy switching mechanisms determined by intelligent LLMs.

256 **Hardware Precision.** Our method benefits from hardware that precisely reaches target poses for  
 257 waypoint tracking and policy switching, leading to high success rates (Appendix A.10). End-to-end  
 258 approaches can absorb hardware inaccuracies through large-scale data collection.

259 **Pose Estimation.** Exotic objects, reflective surfaces, or partial occlusions challenge off-the-shelf  
 260 pose estimation models. However, these models can be fine-tuned relatively easily using only images  
 261 and CAD data, without requiring expensive robot demonstrations.

## 262 References

- 263 [1] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion  
264 policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics*  
265 *Research*, 44(10-11):1684–1704, 2025.
- 266 [2] K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail, M. R. Equi, C. Finn,  
267 N. Fusai, M. Y. Galliker, D. Ghosh, L. Groom, K. Hausman, b. ichter, S. Jakubczak, T. Jones,  
268 L. Ke, D. LeBlanc, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, A. Z. Ren, L. X.  
269 Shi, L. Smith, J. T. Springenberg, K. Stachowicz, J. Tanner, Q. Vuong, H. Walke, A. Walling,  
270 H. Wang, L. Yu, and U. Zhilinsky.  $\pi_{0.5}$ : a vision-language-action model with open-world  
271 generalization. In *Proceedings of The 9th Conference on Robot Learning*, volume 305 of  
272 *Proceedings of Machine Learning Research*, pages 17–40. PMLR, 2025.
- 273 [3] Y. Li, H. Jiang, J. Xia, H. Zhang, J. Du, Y. Zhou, J. Zeng, C. Hao, J. Ren, Q. Yu, et al.  
274 Forcevla2: Unleashing hybrid force-position control with force awareness for contact-rich ma-  
275 nipulation. *arXiv preprint arXiv:2603.15169*, 2026.
- 276 [4] A. O’Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta,  
277 A. Mandlekar, A. Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x mod-  
278 els: Open x-embodiment collaboration 0. In *IEEE International Conference on Robotics and*  
279 *Automation (ICRA)*, pages 6892–6903, 2024.
- 280 [5] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany,  
281 M. K. Srirama, L. Y. Chen, K. Ellis, et al. DROID: A large-scale in-the-wild robot manipulation  
282 dataset. In *Robotics: Science and Systems*, 2024.
- 283 [6] T. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation with  
284 low-cost hardware. *Robotics: Science and Systems (RSS)*, 2023.
- 285 [7] NVIDIA et al. GR00T N1: An open foundation model for generalist humanoid robots. *arXiv*  
286 *preprint arXiv:2503.14734*, 2025.
- 287 [8] N. Carion, L. Gustafson, Y.-T. Hu, S. Debnath, R. Hu, D. Suris, C. Ryali, K. V. Alwala,  
288 H. Khedr, A. Huang, et al. Sam 3: Segment anything with concepts. *arXiv preprint*  
289 *arXiv:2511.16719*, 2025.
- 290 [9] B. Wen, W. Yang, J. Kautz, and S. Birchfield. Foundationpose: Unified 6d pose estimation and  
291 tracking of novel objects. In *Proceedings of the IEEE/CVF conference on computer vision and*  
292 *pattern recognition*, pages 17868–17879, 2024.
- 293 [10] H. Li, Y. Cui, and D. Sadigh. How to train your robots? the impact of demonstration modality  
294 on imitation learning. In *2025 IEEE International Conference on Robotics and Automation*  
295 *(ICRA)*, pages 1113–1120, 2025.
- 296 [11] S. Belkhale, Y. Cui, and D. Sadigh. Data quality in imitation learning. In *Advances in Neural*  
297 *Information Processing Systems*, volume 36, pages 80375–80395. Curran Associates, Inc.,  
298 2023.
- 299 [12] J. Luo, C. Xu, J. Wu, and S. Levine. Precise and dexterous robotic manipulation via human-  
300 in-the-loop reinforcement learning. *Science Robotics*, 10(105), 2025.
- 301 [13] B. Tang, I. Akinola, J. Xu, B. Wen, A. Handa, K. Van Wyk, D. Fox, G. S. Sukhatme, F. Ramos,  
302 and Y. S. Narang. Automate: Specialist and generalist assembly policies over diverse geome-  
303 tries. In *Robotics: Science and Systems*, 2024.
- 304 [14] Y. Tian, J. Xu, Y. Li, J. Luo, S. Sueda, H. Li, K. D. D. Willis, and W. Matusik. Assemble them  
305 all: Physics-based planning for generalizable assembly by disassembly. *ACM Transactions on*  
306 *Graphics*, 41(6):1–11, 2022.

- 307 [15] Y. Tian, K. D. Willis, B. Al Omari, J. Luo, P. Ma, Y. Li, F. Javid, E. Gu, J. Jacob, S. Sueda, et al.  
308 Asap: Automated sequence planning for complex robotic assembly with physical feasibility. In  
309 *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4380–4386.  
310 IEEE, 2024.
- 311 [16] L. Ankile, A. Simeonov, I. Shenfeld, and P. Agrawal. Juicer: Data-efficient imitation learning  
312 for robotic assembly. In *2024 IEEE/RSJ International Conference on Intelligent Robots and  
313 Systems (IROS)*, pages 5096–5103. IEEE, 2024.
- 314 [17] B. Fu, S. K. Leong, X. Lian, and X. Ji. 6d robotic assembly based on rgb-only object pose  
315 estimation. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems  
316 (IROS)*, pages 4736–4742. IEEE, 2022.
- 317 [18] A. Morgan, B. Wen, J. Liang, A. Boularias, A. Dollar, and K. Bekris. Vision-driven compliant  
318 manipulation for reliable; high-precision assembly tasks. In *Robotics: Science and Systems  
319 XVII, RSS2021*. Robotics: Science and Systems Foundation, 2021.
- 320 [19] B. Fu, S. K. Leong, Y. Di, G. Wang, J. Tang, F. Tombari, and X. Ji. Lanpose: Language-  
321 instructed 6d object pose estimation for robotic assembly. In *European Conference on Com-  
322 puter Vision*, pages 43–59. Springer, 2024.
- 323 [20] B. Tang, M. A. Lin, I. A. Akinola, A. Handa, G. S. Sukhatme, F. Ramos, D. Fox, and Y. S.  
324 Narang. IndustReal: Transferring Contact-Rich Assembly Tasks from Simulation to Reality.  
325 In *Proceedings of Robotics: Science and Systems*, 2023.
- 326 [21] G. Schoettler, A. Nair, J. A. Ojea, S. Levine, and E. Solowjow. Meta-reinforcement learning  
327 for robotic industrial insertion tasks. In *2020 IEEE/RSJ International Conference on Intelligent  
328 Robots and Systems (IROS)*, pages 9728–9735. IEEE, 2020.
- 329 [22] T. Bi and R. D’Andrea. Sample-efficient learning to solve a real-world labyrinth game using  
330 data-augmented model-based reinforcement learning. In *2024 IEEE International Conference  
331 on Robotics and Automation (ICRA)*, pages 7455–7460. IEEE, 2024.
- 332 [23] L. Ankile, A. Simeonov, I. Shenfeld, M. Torne, and P. Agrawal. From imitation to refinement-  
333 residual rl for precise assembly. In *2025 IEEE International Conference on Robotics and  
334 Automation (ICRA)*, pages 01–08. IEEE, 2025.
- 335 [24] A. Goyal, V. Blukis, J. Xu, Y. Guo, Y.-W. Chao, and D. Fox. RVT-2: Learning Precise Manip-  
336 ulation from Few Demonstrations. In *Proceedings of Robotics: Science and Systems*, 2024.
- 337 [25] H. L. Brown, G. Hollinger, and S. Lee. Learning hybrid-control policies for high-precision  
338 in-contact manipulation under uncertainty. *arXiv preprint arXiv:2604.19677*, 2026.
- 339 [26] F. Shao, S. Endo, S. Hirche, and F. Ficuciello. Interactive force-impedance control. *IEEE  
340 Robotics and Automation Letters*, 11(5):6488–6495, 2026.
- 341 [27] J. Liang, X. Cheng, and O. Kroemer. Learning preconditions of hybrid force-velocity con-  
342 trollers for contact-rich manipulation. In *Proceedings of The 6th Conference on Robot Learn-  
343 ing*, volume 205 of *Proceedings of Machine Learning Research*, pages 679–689. PMLR, 2023.
- 344 [28] W. Liu, J. Wang, Y. Wang, W. Wang, and C. Lu. Forcemimic: Force-centric imitation learning  
345 with force-motion capture system for contact-rich manipulation. In *2025 IEEE International  
346 Conference on Robotics and Automation (ICRA)*, pages 1105–1112. IEEE, 2025.
- 347 [29] H. Fang, S. Tang, M. Mei, H. Qin, Z. He, J. Chen, Y. Feng, C. Wang, W. Liu, Z. He, et al. Force  
348 policy: Learning hybrid force-position control policy under interaction frame for contact-rich  
349 manipulation. *arXiv preprint arXiv:2602.22088*, 2026.

- 350 [30] M. A. Lee, C. Florensa, J. Tremblay, N. Ratliff, A. Garg, F. Ramos, and D. Fox. Guided  
351 uncertainty-aware policy optimization: Combining learning and model-based strategies for  
352 sample-efficient policy learning. In *2020 IEEE International Conference on Robotics and*  
353 *Automation (ICRA)*, pages 7505–7512, 2020.
- 354 [31] H. Fang, S. Tang, M. Mei, H. Qin, Z. He, J. Chen, Y. Feng, C. Wang, W. Liu, Z. He, et al. Force  
355 policy: Learning hybrid force-position control policy under interaction frame for contact-rich  
356 manipulation. *arXiv preprint arXiv:2602.22088*, 2026.
- 357 [32] L. Wang, J. Zhao, Y. Du, E. H. Adelson, and R. Tedrake. PoCo: Policy composition from and  
358 for heterogeneous robot learning. In *Robotics: Science and Systems*, 2024.
- 359 [33] H. Chen, J. Xu, H. Chen, K. Hong, B. Huang, C. Liu, J. Mao, Y. Li, Y. Du, and K. Driggs-  
360 Campbell. Multi-modal manipulation via multi-modal policy consensus. *arXiv preprint*  
361 *arXiv:2509.23468*, 2025.
- 362 [34] Y. Wang, Y. Zhang, M. Huo, R. Tian, X. Zhang, Y. Xie, C. Xu, P. Ji, W. Zhan, M. Ding, et al.  
363 Sparse diffusion policy: A sparse, reusable, and flexible policy for robot learning. 2024.
- 364 [35] R. Römer, Y. Zhang, Y. Li, and A. P. Schoellig. Clare: Continual learning for vision-language-  
365 action models via autonomous adapter routing and expansion. *IEEE Robotics and Automation*  
366 *Letters*, 2026.
- 367 [36] A. Mandlekar, S. Nasiriany, B. Wen, I. Akinola, Y. Narang, L. Fan, Y. Zhu, and D. Fox. Mim-  
368 icgen: A data generation system for scalable robot learning using human demonstrations. In  
369 *Proceedings of The 7th Conference on Robot Learning*, volume 229 of *Proceedings of Machine*  
370 *Learning Research*, pages 1820–1864. PMLR, 2023.
- 371 [37] A. Mandlekar, C. R. Garrett, D. Xu, and D. Fox. Human-in-the-loop task and motion planning  
372 for imitation learning. In *Proceedings of The 7th Conference on Robot Learning*, volume 229  
373 of *Proceedings of Machine Learning Research*, pages 3030–3060. PMLR, 2023.
- 374 [38] C.-H. Yeh, T.-S. Nan, R. Vuorio, W. Hung, H. Y. Wu, S.-H. Sun, and P.-C. Hsieh. Action-  
375 constrained imitation learning. In *Forty-second International Conference on Machine Learn-*  
376 *ing*, 2025.
- 377 [39] T. Lee, A. Wagenmaker, K. Pertsch, P. Liang, S. Levine, and C. Finn. Roboreward: General-  
378 purpose vision-language reward models for robotics. *arXiv preprint arXiv:2601.00675*, 2026.
- 379 [40] A. Liang, Y. Korkmaz, J. Zhang, M. Hwang, A. Anwar, S. Kaushik, A. Shah, A. S. Huang,  
380 L. Zettlemoyer, D. Fox, et al. Robometer: Scaling general-purpose robotic reward models via  
381 trajectory comparisons. *arXiv preprint arXiv:2603.02115*, 2026.
- 382 [41] P. J. Ball, L. Smith, I. Kostrikov, and S. Levine. Efficient online reinforcement learning with  
383 offline data. In *International Conference on Machine Learning*, pages 1577–1594. PMLR,  
384 2023.
- 385 [42] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy  
386 deep reinforcement learning with a stochastic actor. In *International conference on machine*  
387 *learning*, pages 1861–1870. Pmlr, 2018.
- 388 [43] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*,  
389 2016.
- 390 [44] X. Chen, C. Wang, Z. Zhou, and K. W. Ross. Randomized ensembled double q-learning:  
391 Learning fast without a model. In *International Conference on Learning Representations*,  
392 2021.
- 393 [45] T. Collins and A. Bartoli. Infinitesimal plane-based pose estimation. *International Journal of*  
394 *Computer Vision*, 109(3):252–286, 2014.

- 395 [46] J. Stranghöner, P. Hartmann, M. Braun, S. Wrede, and K. Neumann. Share-rl: Structured,  
396 interactive reinforcement learning for contact-rich industrial assembly tasks. *arXiv preprint*  
397 *arXiv:2509.13949*, 2025.
- 398 [47] X. Zhang, S. Jin, C. Wang, X. Zhu, and M. Tomizuka. Learning insertion primitives with  
399 discrete-continuous hybrid action space for robotic assembly tasks. In *2022 International*  
400 *conference on robotics and automation (ICRA)*, pages 9881–9887. IEEE, 2022.
- 401 [48] M. Caron, H. Touvron, I. Misra, H. Jegou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging  
402 properties in self-supervised vision transformers. In *IEEE/CVF International Conference on*  
403 *Computer Vision (ICCV)*, pages 9630–9640, 2021.
- 404 [49] S. Dasari, O. Mees, S. Zhao, M. K. Srirama, and S. Levine. The ingredients for robotic dif-  
405 fusion transformers. In *2025 IEEE International Conference on Robotics and Automation*  
406 *(ICRA)*, pages 15617–15625. IEEE, 2025.
- 407 [50] Y. Chen, S. Tian, S. Liu, Y. Zhou, H. Li, and D. Zhao. ConRFT: A Reinforced Fine-tuning  
408 Method for VLA Models via Consistency Policy. In *Proceedings of Robotics: Science and*  
409 *Systems*, 2025.
- 410 [51] E. Nikishin, M. Schwarzer, P. D’Oro, P.-L. Bacon, and A. Courville. The primacy bias in  
411 deep reinforcement learning. In *Proceedings of the 39th International Conference on Machine*  
412 *Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 16828–16847.  
413 PMLR, 2022.
- 414 [52] D. Coleman, I. Sucas, S. Chitta, and N. Correll. Reducing the barrier to entry of complex  
415 robotic software: a moveit! case study. *arXiv preprint arXiv:1404.3785*, 2014.
- 416 [53] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In  
417 *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–  
418 5033. IEEE, 2012.
- 419 [54] M. Lepert, J. Fang, and J. Bohg. Phantom: Training robots without robots using only human  
420 videos. In *Proceedings of The 9th Conference on Robot Learning*, volume 305 of *Proceedings*  
421 *of Machine Learning Research*, pages 4545–4565. PMLR, 2025.
- 422 [55] A. Maddukuri, Z. Jiang, L. Y. Chen, S. Nasiriany, Y. Xie, Y. Fang, W. Huang, Z. Wang, Z. Xu,  
423 N. Chernyadev, S. Reed, K. Goldberg, A. Mandlekar, L. Fan, and Y. Zhu. Sim-and-real co-  
424 training: A simple recipe for vision-based robotic manipulation. In *Proceedings of Robotics:*  
425 *Science and Systems (RSS)*, 2025.
- 426 [56] R. Cadene, S. Alibert, A. Soare, Q. Gallouedec, A. Zouitine, S. Palma, P. Kooijmans, M. Ar-  
427 actingi, M. Shukor, D. Aubakirova, M. Russi, F. Capuano, C. Pascal, J. Choghari, J. Moss,  
428 and T. Wolf. Lerobot: State-of-the-art machine learning for real-world robotics in pytorch.  
429 <https://github.com/huggingface/lerobot>, 2024.
- 430 [57] D. S. J. Pro, O. Hausdörfer, R. Römer, M. Dösch, M. Schuck, and A. P. Schoellig. Crisp-  
431 compliant ros2 controllers for learning-based manipulation policies and teleoperation. *IEEE*  
432 *Robotics and Automation Practice*, 2026.

## 433 A Appendix

### 434 A.1 Data collection setups



Figure 5: Data collection setups for our different scenes. Our setup requires a fixed scene for data collection. For Lego, we use the base plate to define the layout. For the other tasks, we 3D print small fixtures to keep the part in place.

### 435 A.2 Policy Switching

436 We can choose different mechanisms to trigger the learned policy. Distance-based thresholds depend on the measurements and noise of the depth sensor and can be unreliable for reflective objects.  
437  
438 Thresholds based on the uncertainty of the pose estimation [30] suffer the same issue, and additionally require calibration of the uncertainty for different objects, as the uncertainty depends on the  
439 object and deployment conditions. For our chosen tasks, we find that a F/T-based threshold as used  
440 in [29, 46, 47] is most reliable:  $\|\mathbf{F}\|_2 \geq \mathcal{F}_{\text{threshold}}$ . Thus, we start the critical segment once contact  
441 is made.  
442

443 For determining the success of the learned policy for the critical segment, we use switching based  
444 on the learned Q-function as described in the main text.

445 **A.3 Safety Filter**

446 We constrain the end-effector position to a sphere with center  $\mathbf{s}_{\text{safe}}$  and radius  $r_{\text{safe}}$  during exploration,  
 447 yielding the finally applied action:

$$\mathbf{a}'_t = \begin{cases} v \cdot \frac{\mathbf{s}_{\text{safe}} - \mathbf{s}_t}{\|\mathbf{s}_{\text{safe}} - \mathbf{s}_t\|_2} & \text{if } \|\mathbf{s}_{\text{safe}} - \mathbf{s}_t\|_2 > r_{\text{safe}} \\ \mathbf{a}_t & \text{otherwise.} \end{cases} \quad (5)$$

448 This formulation redirects the applied action toward the safe region whenever exploration drives the  
 449 end-effector outside the defined sphere. This promotes dense exploration within the critical segment.

450 **A.4 Success Prediction**

451 The task-completion signal stems from the success classifier, which is the Q-function. The Q-  
 452 function predicts what state-action pair results in the sparse success reward. Therefore, we use a  
 453 threshold  $\lambda_{\text{success}}$  to detect successful insertion:

$$Q(\mathbf{s}_t, \pi_\theta(\mathbf{s}_t)) > \lambda_{\text{success}} \quad (6)$$

454 where we average the Q-values across the ensemble of  $N_Q$  Q functions:

$$Q(\mathbf{s}, \mathbf{a}) = \frac{1}{N_Q} \sum_{i=1}^{N_Q} Q_{\theta_i}(\mathbf{s}, \mathbf{a}). \quad (7)$$

455 In practice  $\lambda_{\text{success}}$  needs to be tuned for deployment. If  $\lambda_{\text{success}}$  is set too low the policy stops the  
 456 insertion prematurely, and set too high it fails to identify successes. To avoid the first case (false  
 457 positive), we choose  $\lambda_{\text{success}}$  high. To avoid the latter case (false negative), we additionally use a  
 458 Hysteresis-style maximum detector with two thresholds  $\lambda_{\text{high}}$  and  $\lambda_{\text{low}}$ :

$$\mathcal{H}(t) := \exists t' < t : Q(\mathbf{s}_{t'}, \pi_\theta(\mathbf{s}_{t'})) > \lambda_{\text{high}} \wedge Q(\mathbf{s}_t, \pi_\theta(\mathbf{s}_t)) < \lambda_{\text{low}}, \quad (8)$$

459 where we use  $\lambda_{\text{high}} = 0.8R_{\text{Success}}$  and  $\lambda_{\text{low}} = 0.6R_{\text{Success}}$ , with  $R_{\text{Success}}$  the sparse success reward.  
 460 This term works, since we use a sparse reward setting and the Q-function decreases sharply after  
 461 reaching its maximum value at the insertion point. Overall, we get the following classifier:

$$\text{success} = \begin{cases} 1 & \text{if } Q(\mathbf{s}_t, \pi_\theta(\mathbf{s}_t)) > \lambda_{\text{success}} \\ 1 & \text{if } \mathcal{H}(t) \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

462 Note that this approach has a significant advantage over other methods that require a separately  
 463 trained reward classifier prior to policy training. Instead of manually recording successful and failed  
 464 episodes before training starts, our success classifier is a natural byproduct of policy training.

465 **Results.** We show the resulting Q-value predictions  $Q_\phi(\mathbf{s}_t, \pi_\theta(\mathbf{s}_t))$  as a function of policy timesteps  
 466 for truncated and successful trajectories in the training data in Figure 6. The separation between  
 467 successful and truncated trajectories is clearly visible, motivating our method of success classifica-  
 468 tion.

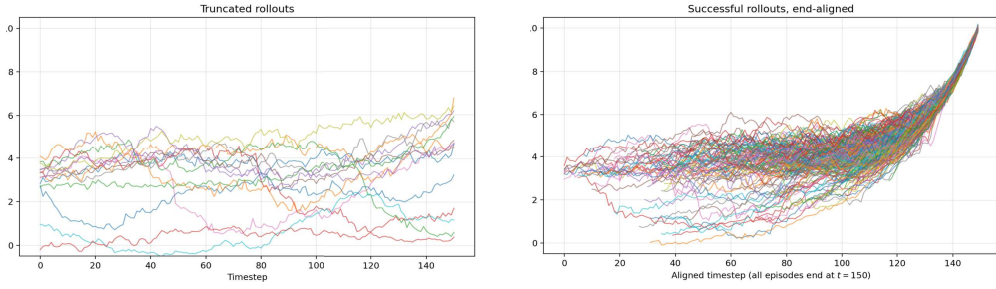


Figure 6: Q-function predictions for truncated and successful rollouts. Successful episodes are aligned such that successful termination occurs at timestep  $t = 150$ . The maximum Q-Value is determined by the sparse reward  $R_{\text{success}} = 10$ .

#### 469 A.5 Evaluation of Pose Estimation

470 We evaluate the pose estimation (FoundationPose [9]) quantitatively on the fan cover task. For  
 471 the evaluation protocol, we placed the fan cover at a fixed position in the scene. We then run pose  
 472 estimation from a distance of  $\approx 30$  cm (coarse) and  $\approx 10$  cm (fine) each 50 times from different end-  
 473 effector positions. Quantitative results are reported in Figure 7 and qualitative results in Figure 8.  
 474 For the fine estimates, all estimates fall within a range of  $\pm 1$  mm and  $\pm 3.3^\circ$ .

475 **Critical segment:** The critical segment is determined by the uncertainty in the pose estimation  
 476 — effectively, we correct for this uncertainty using robot data. In practice, we choose the critical  
 477 segment to be larger than the expected uncertainty.

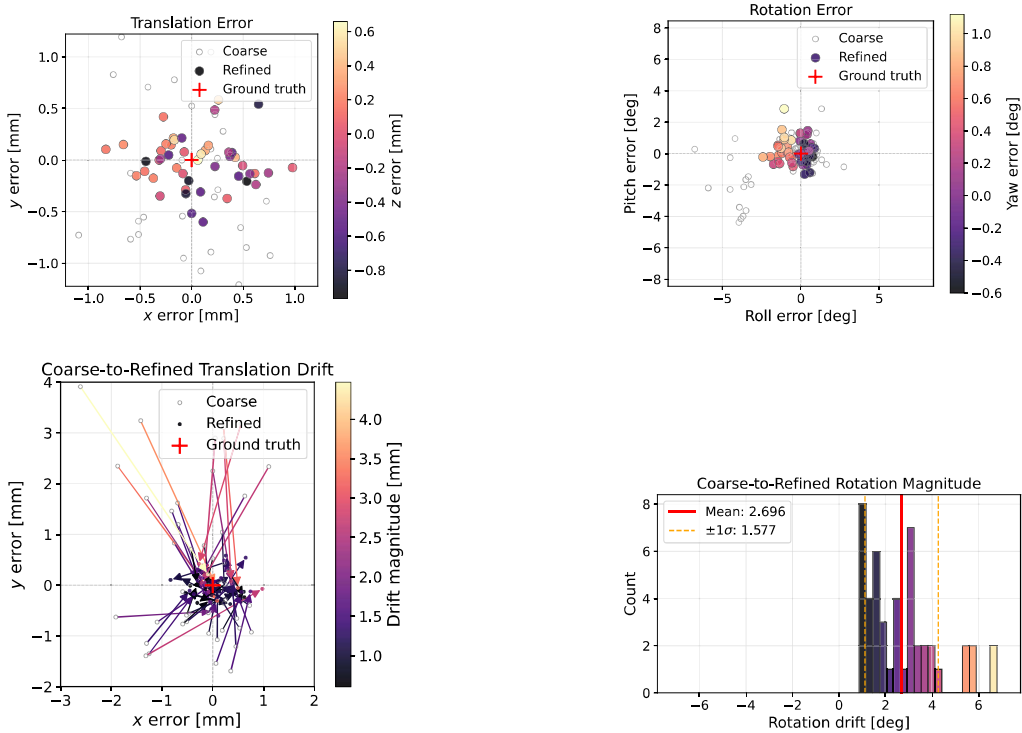


Figure 7: Quantitative pose estimation evaluation. **Top row:** Pose estimation distribution for coarse and fine estimates. **Bottom row:** Correction of pose errors between coarse and fine estimates for positions and orientations.

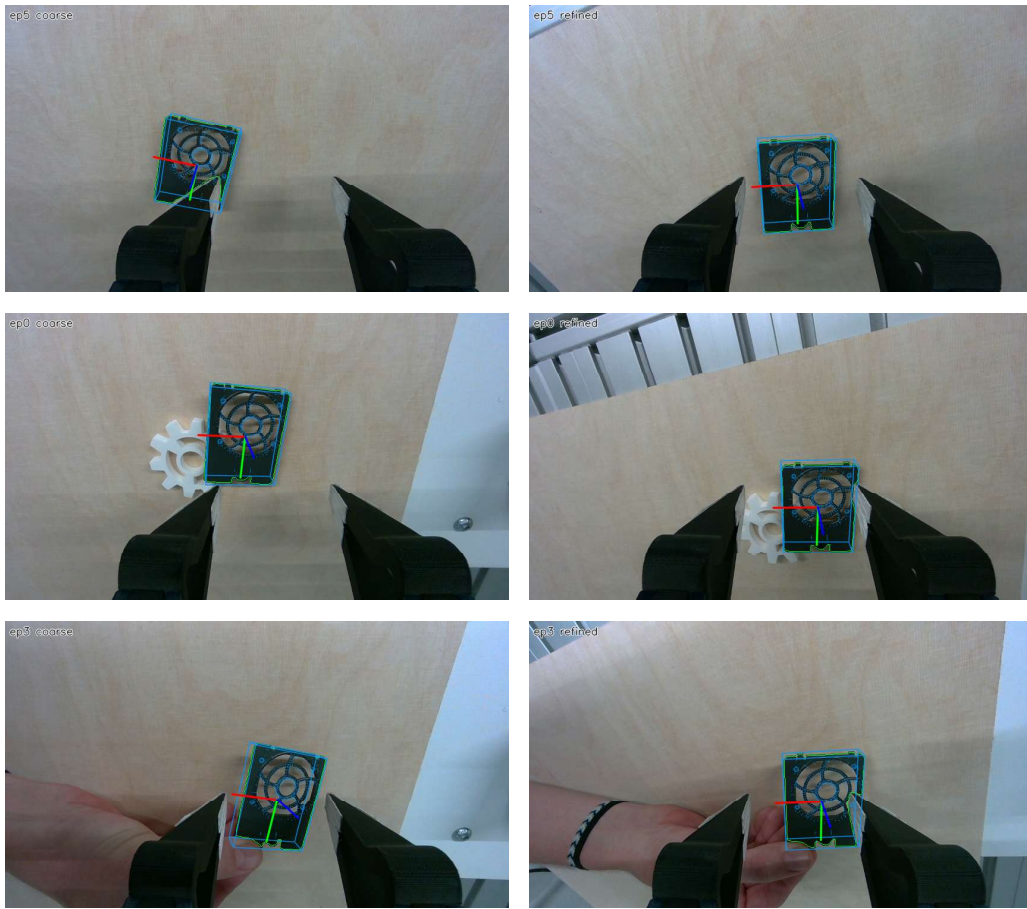


Figure 8: Qualitative pose estimation results including partial occlusions and arbitrary poses. **Left:** Coarse estimates from larger distance. **Right:** Fine estimates from closer distance.

478 **A.6 Out-of-distribution (OOD) evaluations**

479 We show the out-of-distribution settings for quantitative evaluations in Figure 9, and additional  
480 qualitative scenarios only for our method in Figure 10.

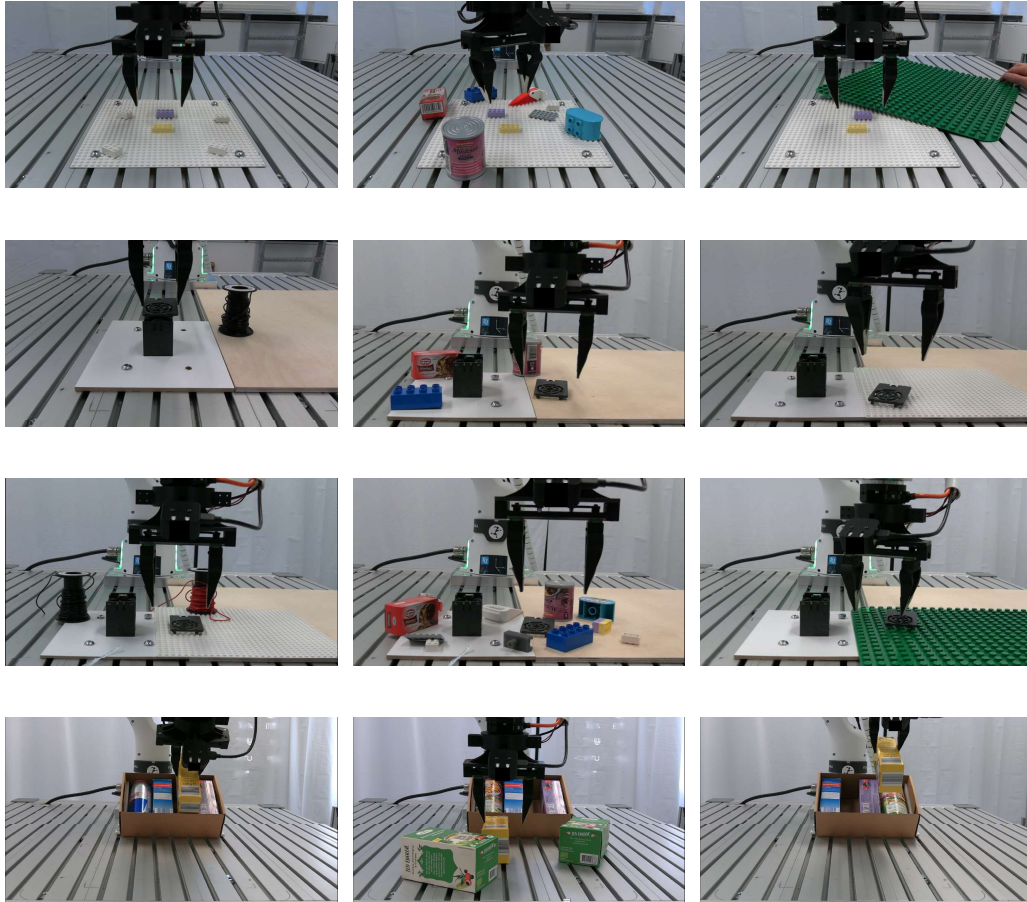


Figure 9: Overview of our out-of-distribution settings used for the quantitative evaluations. For each task, we include evaluations with slight scene distractors, severe distractors, and workspace color changes. **Tasks from top row to bottom:** Lego Stacking, fan cover, fan cover (difficult), shelf stocking. The out-of-distribution settings for the *fan cover* tasks were included in the data collection (in-distribution) for the *fan cover (difficult)* task to test how well training on distractors generalizes.

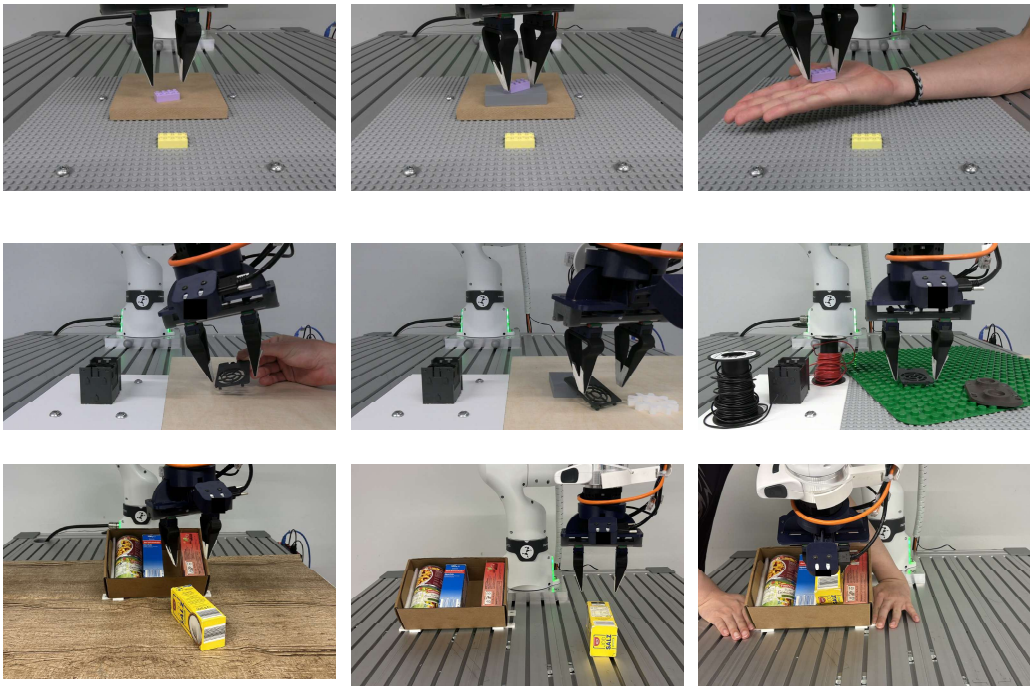


Figure 10: Further qualitative evaluations only performed with method. We place the parts at different heights, perform dynamic grasping from a human hand, and fully randomize the object poses in 6D – all without additional data collection compared to the in-distribution experiments.

481 **A.7 Experimental Details**

482 For all in-distribution evaluations, we place the base object at a fixed position in the workspace.  
 483 For our method, we still use pose estimation to estimate the base object to demonstrate the full  
 484 pipeline, except for the fan cover task, where we assume the base is fixed. The inserted object is  
 485 randomized for every rollout in a given region and is always estimated for our method. For the  
 486 quantitative out-of-distribution evaluations, the same applies. Here, the positions of the distractor  
 487 objects are randomized for each rollout. All learned policies are actuated in four dimensions (three  
 488 translations and one rotation). This made teleoperation for precise tasks significantly easier, and  
 489 for fair comparison, we applied the same procedure to all baselines and our method. However,  
 490 for our method, we also show full 6D manipulation trajectories in the additional out-of-distribution  
 491 evaluations.

492 For all baselines, we use their LeRobot implementations [56]. The data used for imitation learning is  
 493 filtered for zero-action frames. We tried to replicate HIL-SERL’s success rates for our tasks, which  
 494 was difficult. We simplified the shelf stacking task slightly by assuming the yellow salt package was  
 495 pregrasped.

496 **A.8 Training parameters**

497 Training and architectural details for our method are provided in Table 4.

498 **Training.** We train our policies, DP, and DiTFlow on a local workstation with an NVIDIA RTX5090  
 499 GPU. Our policy trains roughly 1 h, while the others train 5 h.  $\pi_{0.5}$  is full-finetuned on two cluster  
 500 nodes with four H100 GPUs each for 18 h. The equivalent cost for a commercial cloud compute  
 501 provider would roughly be \$200 as at the time of writing.

502 **Hardware.** Our hardware is a FRANKA RESEARCH 3. The force-torque sensor we use is the  
 503 BOTA-DENS-IND2-B4. The wrist-mounted camera is an Intel Realsense D405.

Parameter	Value	Explanation
$p$	0.8	Probability of choosing a random action (during data collection).
$L_{\text{trunc}}$	150	Episode truncation length (steps).
$R_{\text{success}}$	10	Task completion reward.
Batch size	512	Samples per optimization step.
Q learning rate	$5 \times 10^{-4}$	Adam step size for critic ensemble.
Policy learning rate	$3 \times 10^{-4}$	Adam step size for actor.
$\alpha$	$1 \times 10^{-3}$	Fixed entropy regularization weight.
$\gamma$	0.97	Reward discount factor.
$\tau$	0.999	Polyak averaging coefficient.
UTD	5	Update-to-data ratio, equals number of epochs.
$N_Q$	10	Critic ensemble size.
$N_{Q,S}$	2	Critic ensemble subset size.
Actor	MLP	Two-layers, hidden size 128, ReLU.
Q-functions	MLP	Ensemble of $N_Q = 10$ , two-layers hidden size 128, LayerNorm + ReLU.
$f_{\text{lowlevel}}$	1000Hz	CRISP controllers [57] for direct torque-control.
$f_{\text{policy}}$	15Hz	Control frequency of the DRL policy.

Table 4: Data collection and training hyperparameters.

504 **A.9 Camera Cropping**

505 We show the visual inputs received by our policies in Figure 11. The used image resolution is  
506  $224 \times 224$ .



Figure 11: Visual inputs used for our policies. Policies for our method receive a single viewpoint from a wrist-mounted camera as input that is cropped to the critical insertion region of the task.

507 **A.10 Failure cases**

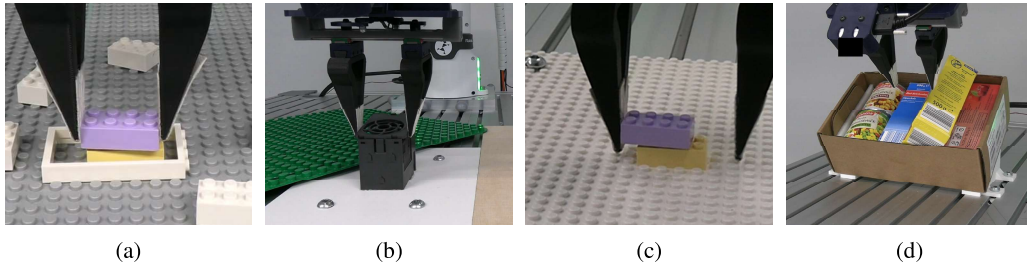


Figure 12: Policy failure cases for our method that partially concern baselines as well. **(a) + (b)** Our policy fails if distractors are placed very close to the object within the cropped region of the wrist image taking the policy OOD. **(b)** Hardware limitations, such as objects sticking to the gripper, play a role when success rates reach 100%. **(c)** In the one failure case of our method for the shelf insertion, the policy went OOD during insertion, possibly because objects in the shelf moved.

508 **A.11 Data Collection Ablations (Sensor Modalities)**

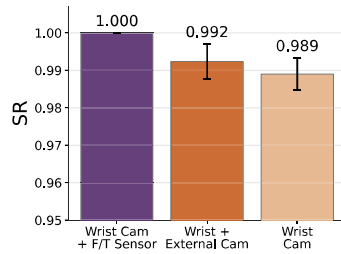


Figure 13: Policy success rate for different combinations of sensing modalities.

509 **A.12 Episodic maximum torques**

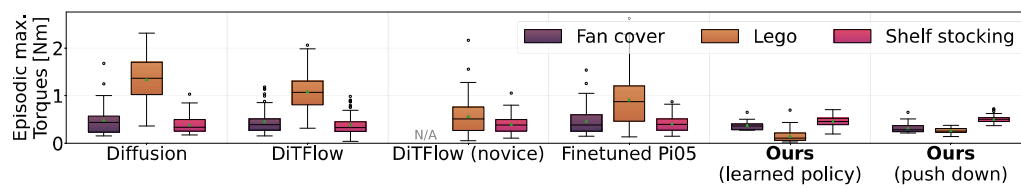


Figure 14: Maximum torques applied during policy rollouts. **Ours** applies significantly less torques to the objects than baselines. In the figure, we divide ours between the learned insertion policy and the planned push-down motion.