

---

# DexFlyWheel: A Scalable and Self-improving Data Generation Framework for Dexterous Manipulation

---

Anonymous Author(s)

Affiliation

Address

email



Figure 1: DexFlyWheel generates overall 2000 demonstrations from only one human demonstration per task, including various scenarios. The policies trained on our dataset are successfully deployed in the real-world.

## Abstract

1 Dexterous manipulation is critical to advancing robot capabilities in real-world  
2 applications, yet diverse and high-quality datasets remain scarce. Existing data  
3 collection methods either rely on human teleoperation or require significant human  
4 engineering, or are merely limited to grasping, restricting their scalability and gen-  
5 eralization. In this paper, we introduce DexFlyWheel, a scalable data generation  
6 framework that employs a self-improving cycle to iteratively expand data diversity.  
7 Starting from efficient seed demonstrations warmup, our framework expands data  
8 diversity via multiple iterations in the self-improving cycle. Each iteration follows  
9 a closed-loop pipeline that combines imitation learning, reinforcement learning,  
10 rollout trajectory collection, and data augmentation. At each iteration, we first use  
11 imitation learning to extract behavioral priors from demonstrations and employ  
12 reinforcement learning to enhance generalization. Based on our policy, we rollout  
13 trajectories in simulation and then augment these across different environments and  
14 objects positions. As iterations progress, our framework generates more diverse  
15 data, including various objects, environments, and object positions. Experi-  
16 mental results show that policies trained on our dataset achieve an average success  
17 rate of 81.9% on the challenge test sets, with a real-world transfer success rate  
18 of 78.3% on dual-arm lift tasks. Videos can be found on our project website  
19 <https://DexFlyWheel.github.io>.

20 **1 Introduction**

21 Learning from Demonstration (LfD) [1] has become increasingly prevalent in robotics. This approach  
22 typically requires collecting human demonstration data through teleoperation, where humans directly  
23 control robots to perform tasks. In dexterous manipulation particularly [2, 3, 4], the higher degrees  
24 of freedom and richer contact interactions demand substantially more data. Recent work has shown  
25 that training large models with extensive datasets can achieve more challenging tasks and better  
26 generalization [5, 6, 7, 8]. However, collecting sufficiently large and diverse demonstration datasets  
27 faces several critical limitations: it requires significant human resources, suffers from weak cross-  
28 embodiment transfer (where changes in robot embodiment or sensor placement can invalidate  
29 collected data), and typically constrains data collection to laboratory settings. While some researchers  
30 have explored using portable motion capture devices for in-the-wild data collection [9], this approach  
31 still requires substantial human involvement and remains limited by cross-embodiment gap. Recently,  
32 simulation has emerged as a promising solution to address data challenges in robotics [10, 11, 12, 13,  
33 14, 15]. It offers numerous advantages: parallel data collection at scale, easy modification of robot  
34 embodiments and sensor configurations, and domain randomization for diverse data generation. Some  
35 researchers have employed optimization algorithms or heuristic planning methods [16] to generate  
36 demonstrations in simulation. However, these approaches often produce lower-quality trajectories  
37 compared to human demonstrations and face the limitations in significant optimization challenges,  
38 making them inadequate for dynamic dexterous tasks. While reinforcement learning (RL) has been  
39 widely adopted for training policies and collecting data in simulation [14, 17, 18, 19, 20], purely  
40 RL-trained models often exhibit non-human-like behaviors, leading to less robust manipulation and  
41 increased sim-to-real transfer challenges. Moreover, RL faces exploration difficulties and relies  
42 heavily on reward engineering challenges that are particularly acute in dexterous manipulation. In the  
43 context of high-quality dexterous manipulation data generation, these limitations make pure RL-based  
44 methods less suitable for generating high-quality data.

45 Given these simulation challenges, researchers have begun exploring the teleoperation with replay-  
46 mechanism [10, 21], where humans teleoperate simulated robots to collect training data and then use  
47 spatial transformations to synthesize new trajectories. Although this approach captures relatively  
48 high-quality data while leveraging simulation-based data augmentation, it typically allows only simple  
49 trajectory editing or domain randomization, which constrains the diversity of generated trajectories,  
50 especially when handling different objects. These limitations fundamentally restrict the collected  
51 data within the scope of collected human demonstrations, motivating us to rethink the role of human  
52 demonstrations in data generation pipelines. We observe that manipulating different objects only  
53 induces minor changes in the manipulation trajectories. This suggests that human demonstrations  
54 could serve as a strong behavioral prior, while RL could refine these priors to handle unseen scenarios  
55 through targeted exploration.

56 Building on this insight, we propose DexFlyWheel, a framework that integrates imitation learning  
57 (IL) with residual reinforcement learning (RL). IL captures human-like behaviors from demonstra-  
58 tions, while RL enables generalization by learning adaptive adjustments, enabling the generation of  
59 human-like and environmentally robust trajectories. We further enhance data diversity through data  
60 augmentation across spatial and environments. By embedding them in a self-improving cycle, we  
61 achieve continuous expansion of data diversity and enhancement of policy performance. Our main  
62 contributions include:

- 63 • We present DexFlyWheel, a novel framework that efficiently generates diverse dexterous manipula-  
64 tion data through iterative self-improving. Starting from a single human demonstration per task,  
65 our framework efficiently generates diversity and high-quality demonstrations while maintaining  
66 human-like behavior patterns.
- 67 • We demonstrate the effectiveness of our framework across multiple dexterous manipulation tasks.  
68 We generate over 2,000+ successful demonstrations from just one human demonstration per task  
69 with 500+ diverse scenarios, including 80 objects, 12 environments and 15 spatial configurations.  
70 Our approach significantly outperforms baseline data generation methods in both diversity and  
71 robustness.
- 72 • We validate that policies trained with our generated data achieve over average 81.9% success rate  
73 in challenging test sets, significantly outperforming policies trained with data from alternative  
74 baselines. This superior performance demonstrates the framework’s ability to enhance policy  
75 generalization across diverse scenarios.

76 **2 Related Work**

77 **2.1 Dexterous Manipulation**

78 Dexterous manipulation with multi-fingered robotic hands remains a significant challenge in  
79 robotics [22, 23, 24, 25, 26], largely constrained by high-quality demonstration data scarcity. While  
80 the prevailing approach employs reinforcement learning (RL) to develop manipulation skills, this  
81 method frequently encounters efficiency limitations and exploration challenges [17, 27, 28]. Re-  
82 searchers have explored human video demonstrations [29, 30, 31, 32, 33, 34, 35, 36], but morpho-  
83 logical differences between human and robotic hands create substantial transfer barriers. Human  
84 teleoperation has emerged as a promising alternative for collecting expert trajectories for imitation  
85 learning (IL) [37, 21, 38, 10, 39, 40], effectively capturing expert actions in native robot morphology.  
86 Nevertheless, existing approaches still struggle with data collection efficiency or require extensive  
87 reward engineering, emphasizing the need for high-quality dexterous manipulation datasets.

88 **2.2 Robotic Data Generation in Simulation**

89 Current approaches for collecting robotic demonstrations in simulation face significant limitations  
90 when applied to dexterous manipulation. Motion planning-based methods, while effective for gripper-  
91 based systems [16, 41, 42, 11, 18], struggle with the high-dimensional action space and complex  
92 contact dynamics of multi-fingered manipulation. Large language models (LLMs) [15, 43, 44, 45] can  
93 effectively generate high-level command, but they demonstrate significant limitations when confronted  
94 with high-degree-of-freedom dexterous hands, unable to provide the fine-grained guidance necessary  
95 for coordinated finger-level control. Reinforcement learning-based approaches generate data through  
96 environmental interaction but are severely hampered by scaling challenges inherent to dexterous  
97 manipulation, where the high degrees of freedom in multi-fingered hands create exponentially  
98 large action spaces that make effective training extremely difficult [14, 42, 46]. Replay-based  
99 methods [21, 10, 47] attempt to edit existing demonstrations to new scenarios but face fundamental  
100 scalability constraints, as they merely implement spatial transformations of recorded trajectories  
101 without the ability to explore novel manipulation strategies beyond the original demonstrations.  
102 These collective limitations underscore the critical need for more efficient and scalable methods to  
103 generate diverse, high-quality dexterous manipulation data in simulation. Our proposed DexFlyWheel  
104 framework addresses these challenges by integrating the complementary strengths of these approaches  
105 within a cyclical policy learning paradigm, providing a more effective solution for dexterous skill  
106 acquisition that overcomes the individual shortcomings of existing methods.

107 **3 Task Formulation**

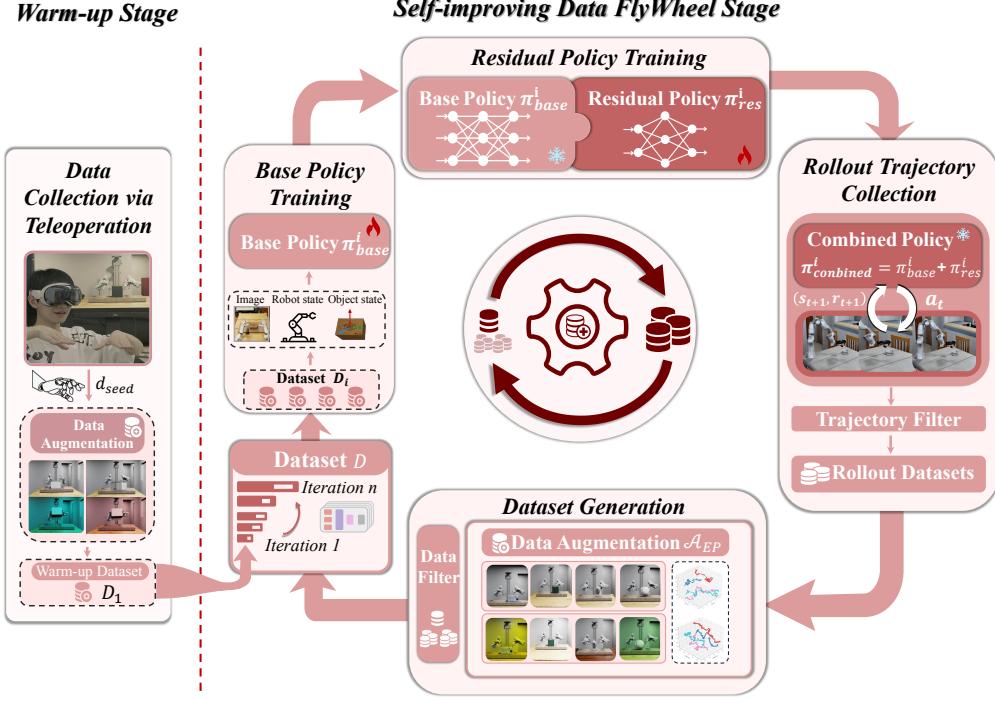
108 To address the challenge of generating high-quality synthetic data for robotic manipulation tasks, we  
109 train policy models for each manipulation task in simulation environments and use these policies to  
110 collect demonstrations. We formulate each manipulation task as a Markov Decision Process (MDP)  
111  $\mathcal{M} = (\mathbf{S}, \mathbf{A}, \pi, \mathcal{T}, R, \gamma, \rho, G)$ , where  $\mathbf{S}$  is the state space,  $\mathbf{A}$  is the action space,  $\pi$  is the agent's policy,  
112  $\mathcal{T}(s_{t+1}|s_t, a_t)$  is the transition distribution,  $R$  is the reward function,  $\gamma$  is the discount factor, and  $\rho$  is  
113 the initial state distribution. The policy  $\pi$  conditions on the current state  $s_t$ , and generates robot action  
114 distributions  $a_t$  to maximize the likelihood between the future object states  $(s_{t+1}, s_{t+2}, \dots, s_{t+T})$ .

115 **4 Method**

116 In this section, we present DexFlyWheel, a framework for autonomous data generation that bootstraps  
117 dexterous manipulation through a self-improving policy learning cycle. We begin with an overview  
118 of the DexFlyWheel architecture (Section 4.1) and then detail the two-stage data generation pipeline  
119 (Sections 4.2 and 4.3). Together, these components enable scalable collection of diverse and high-  
120 quality demonstrations.

121 **4.1 Overview**

122 DexFlyWheel aims to generate diverse and high-quality demonstrations across various scenario  
123 configurations, providing broad coverage of objects, environments, and spatial variations, while only  
124 starting with minimal human demonstrations. Specifically, DexFlyWheel consists of two stages:  
125 the warm-up stage and the data FlyWheel stage. Figure 2 illustrates the overall architecture of our  
126 DexFlyWheel framework.



**Figure 2: DexFlyWheel Framework Overview.** Our method consists of two stages: a warm-up stage (left) and a self-improving data flywheel stage (right). In the warm-up stage, seed demonstrations from VR teleoperation are processed through data augmentation to create the initial dataset  $\mathcal{D}_1$ . The data flywheel stage operates as a closed-loop system with four key components: base policy training via imitation learning, residual policy training using reinforcement learning to enhance object generalization, rollout trajectory collection across various object configurations, and data augmentation to incorporate environment and spatial variations. As the flywheel iterates, both data diversity and policy performance progressively improve, enabling comprehensive generalization across objects, environments, and spatial configurations.

127 **Warm-up stage.** A single human demonstration  $d_{seed}$  is collected via a VR-based teleoperation  
 128 system. This seed demonstration is then processed using a multi-dimensional data augmentation  
 129 module  $\mathcal{A}_{EP}$ . Given the human demonstration  $d_{seed}$ , the augmentor generates new demonstrations  
 130 with diverse environment and spatial variations to produce the initial dataset  $\mathcal{D}_1$ .

131 **Self-improving Data FlyWheel stage.** We design a data FlyWheel mechanism to progressively  
 132 enhance both data diversity and policy performance. This stage comprises multiple iterations  
 133  $i = \{1, 2, \dots, n - 1\}$ , at each iteration  $i$ , the following steps are executed: (1) An imitation learning  
 134 policy  $\pi_{base}^i$  is trained on dataset  $\mathcal{D}_i$  (with  $\mathcal{D}_1$  used when  $i = 1$ ). (2) To improve generalization  
 135 to novel objects, a residual reinforcement learning policy  $\pi_{res}^i$  is trained on top of the frozen  $\pi_{base}^i$ ,  
 136 yielding a combined policy  $\pi_{combined}^i = \pi_{base}^i + \pi_{res}^i$ . (3) The combined policy is deployed in simulation  
 137 to generate demonstrations under various object configurations, forming a high-quality rollout dataset  
 138  $\mathcal{D}_O^i$ . (4) Finally,  $\mathcal{D}_O^i$  is further augmented by  $\mathcal{A}_{EP}$  with environment and spatial variations to produce  
 139 the dataset  $\mathcal{D}_{i+1}$  used in the next iteration.

## 140 4.2 Warm-up Stage

141 The first stage aims to generate an initial dataset  $\mathcal{D}_1$  via data augmentation module  $\mathcal{A}_{EP}$ , starting from  
 142 a single human demonstration  $d_{seed}$ . This warm-up stage including two operations:

143 **Data Collection via Teleoperation:** To bootstrap the framework with high-quality seed data, we de-  
 144 sign a VR-based teleoperation system—implemented in simulation using Apple Vision Pro [39]—to  
 145 accurately track human hand, wrist, and head poses. Since large-scale data collection requires heavily  
 146 human effort, we only need a single demonstration, denoted as  $d_{seed}$ . This demonstration serves as

147 the sole seed for subsequent data generation. This makes our method highly efficient in terms of  
 148 human resources while maintaining the quality and diversity of the generated data.

149 **Data Augmentation:** To scale and diversify our dataset, we introduce data augmentation module  $\mathcal{A}_{EP}$ ,  
 150 which builds upon the MimicGen framework [21] and extends it to support multi-dimensional data  
 151 augmentation across various environments and spatial configurations. Data augmentation module  $\mathcal{A}_{EP}$   
 152 takes a base demonstration dataset  $\mathcal{D}$  and augmented scenario configurations  $\mathcal{C}_{aug}$  as input, and outputs  
 153 the augmented dataset  $\mathcal{D}'$ . In the warmup phase, this process is applied to the seed demonstration:  
 154  $\mathcal{D}_1 = \mathcal{A}_{EP}(d_{seed}; \mathcal{C}_1)$ . Through this warmup phase, we establish a foundation of diverse manipulation  
 155 datasets that serves as the starting point for our iterative data FlyWheel mechanism.

### 156 4.3 Self-improving Data FlyWheel Stage

157 The second stage implements a closed-loop data FlyWheel mechanism that iteratively expands  
 158 data diversity across objects, environments and spatial generalization dimensions. At each iteration  
 159  $i \in \{1, 2, \dots, n - 1\}$ , the FlyWheel performs four key operations:

160 **Base Policy Training.** Given the dataset  $\mathcal{D}_i$  from the previous iteration, we employ diffusion-based  
 161 policy as the base policy  $\pi_{base}$  to learning dexterous manipulation skill, obtaining a strong base policy  
 162 for subsequent modules. At each step  $t$ , the base policy  $\pi_{base}^i$  takes the state  $s_t = \{s_t^{vis}, s_t^{obj}, s_t^{prop}\}$   
 163 as inputs, where  $s_t^{vis}$  represents visual input from camera,  $s_t^{obj}$  contains object state information  
 164 including 6D pose (position and orientation) and velocities, and  $s_t^{prop}$  includes robot proprioception  
 165 data consisting of joint positions, velocities, and end-effector poses. The policy outputs a sequence  
 166 of robots actions  $(a_t, a_{t+1}, \dots, a_{t+H})$ , where  $H$  represents the prediction horizon, and each action  
 167  $a_t$  consists of the end-effector 6D pose and target joint angles. Implementation details of the policy  
 168 parameters are provided in Appendix A.1.

169 **Residual Policy Training.** Generalizing to novel objects remains a key challenge in imitation-based  
 170 robotic manipulation, especially when training with limited data. We observe that manipulating  
 171 different objects induces only small changes in the manipulation trajectories, suggesting that a well-  
 172 initialized policy can only require fine-grained adjustments to adapt to new objects. Based on this  
 173 observation, we propose a residual reinforcement learning framework that builds upon the base policy.  
 174 Specifically, we train a residual policy  $\pi_{res}^i$  that takes object state  $s_t^{obj}$  and robot proprioception  $s_t^{prop}$   
 175 as inputs and generates correction actions  $\Delta a = (\Delta a_t, \Delta a_{t+1}, \dots, \Delta a_{t+H})$ . These corrections are  
 176 applied with a scaling factor  $\alpha$  to the trajectory generated by the base policy, forming a combined  
 177 policy:  $\pi_{combined}^i = \pi_{base}^i + \alpha \cdot \pi_{res}^i$ , where  $\tilde{a}_t = a_t + \alpha \cdot \Delta a_t$  represents the action output from the  
 178 combined policy  $\pi_{combined}^i$  at each timestep in the horizon. This approach allows the residual policy  
 179 to start from a reasonable robot actions from  $\pi_{base}$  and focus on learning the fine-grained refinements to  
 180 generalize objects. The reward function  $R$  used to train  $\pi_{res}$  is detailed in Appendix A.2. To make the  
 181 exploration more controllable, we utilize the progressive exploration schedule as proposed by [48].  
 182 Specifically, we define the combined policy during training as:

$$\pi_{combined}(s) = \begin{cases} \pi_{base}(s) + \alpha \cdot \pi_{res}(s) & \text{with probability } \epsilon \\ \pi_{base}(s) & \text{with probability } 1 - \epsilon \end{cases} \quad (1)$$

183 where  $\epsilon$  serves as a mixing coefficient that linearly increases from 0 to 1 over  $T$  training steps,  
 184 progressively shifting control from the base policy to the residual policy. This scheduling is designed  
 185 to stabilize learning by relying on the base policy in early training and gradually incorporating the  
 186 residual corrections as  $\pi_{res}$  becomes more reliable.

187 **Rollout Trajectory Collection.** In robotic manipulation, generating diverse demonstration data  
 188 across different objects represents the most significant challenge for scaling data collection [10].  
 189 Trajectory editing methods cannot adapt to different objects due to their geometry-unaware [10], we  
 190 use the combined policy  $\pi_{combined}$  to rollout in simulation with objects randomization:  $\mathcal{D}_O^i = \{d_j =$   
 191  $\{(s_t, a_t)\}_{t=0}^{T-1} | d_j \sim \pi_{combined}^i\}_{j=1}^K$ , where we collect  $K$  high-quality trajectories by filtering based on  
 192 task success.

193 **Data Augmentation.** In this module, we employ the previously introduced data augmentation module  
 194  $\mathcal{A}_{EP}$  to efficiently augment data in various environment and spatial configurations. Taking the dataset  
 195  $\mathcal{D}_O^i$  and augmented scenario configurations  $\mathcal{C}_{aug}^{i+1}$  as input, the module produces an expanded dataset:  
 196  $\mathcal{D}_{i+1} = \mathcal{A}(\mathcal{D}_O^i; \mathcal{C}_{aug}^{i+1})$ . The expanded dataset  $\mathcal{D}_{i+1}$  is used to train an improved imitation learning  
 197 policy, which serves as the foundation for the next iteration of the FlyWheel.

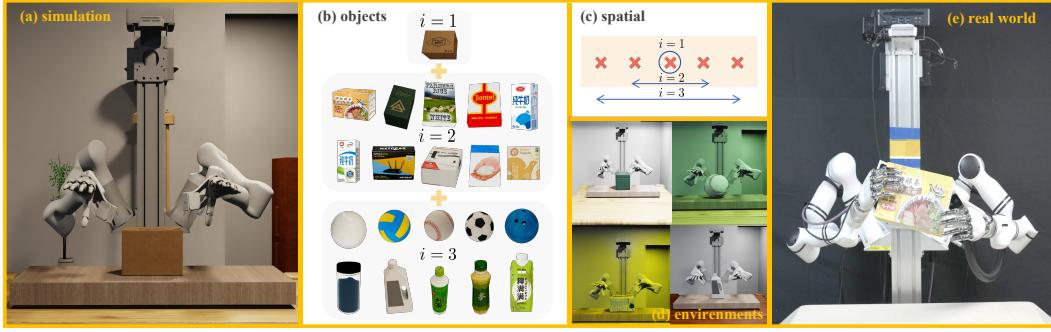


Figure 3: **Experiment Setup.** Taking the dual-arm robot system as an example, **(a)** Our simulation environment. **(b)** Object diversity expansion across iterations, progressing from a single object ( $i=1$ ) to geometrically similar objects ( $i=2$ ) and diverse geometries and physical properties objects ( $i=3$ ). **(c)** Spatial diversity, showing the spatial arrangements. **(d)** Environment diversity, including variations in lighting conditions and tabletop appearances. **(e)** Real-world environment.

## 198 5 Experiments

199 The experiments section are desgin to answer the following questions: **Q1:** How DexFlyWheel  
200 effectively expand data diversity and enhance policy generalization? **Q2:** How does DexFlyWheel  
201 data generation compare to alternatives? **Q3:**What are the strengths of our method, and how does  
202 our system design quantitatively contribute to overall performance? **Q4:** How DexFlyWheel enables  
203 real-world deployment on bimanual dexterous robot system?

### 204 5.1 Experimental Setup

205 **Tasks and Robots.** We evaluate our framework on four dexterous manipulation tasks that assess  
206 both fine-grained control in single-arm settings and coordinated motion in dual-arm scenarios: (1)  
207 Grasp: The robot must grasp the target object and lift it to a height greater than 0.2 meters from the  
208 tabletop. (2) Pour: The robot manipulates a source container to transfer its contents into a target  
209 container, requiring controlled pouring of the contained objects from one receptacle to another. (3)  
210 Lift: The robot performs collaborative manipulation using its two arms to synchronously lift a object  
211 to a minimum height of 15 cm. (4) Handover: The robot performs an intra-agent handover by  
212 transferring an object from one hand to the other through a stable, coordinated motion.

213 For the single-arm Grasp and Pour tasks, we use a Franka Emika Panda robot arm equipped with an  
214 Inspire robotic hand. For the dual-arm Lift and Handover tasks, we use a 7-DoF Real-Man RM75-6F  
215 arm paired with a 6-DoF PsiBot G0-R hand. The robot is equipped with a RealSense D435 camera  
216 mounted on its head, which provides a first-person perspective.

217 **Data Collection and Environment.** For each manipulation task, we collected only a single demon-  
218 stration trajectory using VR teleoperation as our minimal seed data. To ensure the generation of  
219 high-quality data, we employed OmniGibson [49] as our simulation platform, leveraging its realistic  
220 rendering to generates high-quality data. We prepared 80 distinct objects across various categories  
221 and 12 different environments with varying lighting conditions, tabletop appearances. we set the  
222 number of iterations to  $i = \{1, 2, 3\}$ . For each task, we generated 20, 100, and 500 trajectories in the  
223 three iterations, respectively. Simuation setup is visualized in Figure 3. More detailed data collection  
224 and environment setup are provided in Appendix A.3

225 **Evaluation Design** We evaluate our method using two criteria: (1) **data diversity**, the number  
226 of object variations  $O$ , environment variations  $E$  and spatial variations  $P$  in our generated dataset  
227  $D_i$  in each iteration, and the total number of possible scenario configurations ( $O \times E \times P$ ) these  
228 combinations can form from our pipeline; (2) **generalization performance**, the Success Rate ( $SR$ )  
229 of task execution when policies trained on our generated datasets  $D_i$ . It is calculated as the ratio  
230 of successful task completion to the total attempts. A better policy will lead to a higher Success  
231 Rate. Specifically, we construct two types of test sets. First, the multi-factor generalization test set  
232 ( $T_{OEP}$ ) contains 40 unseen scenario configurations that simultaneously incorporate all three types of

233 variations: object, environment, and spatial arrangements. Second, the object generalization test set  
234 ( $T_O^i$ ) evaluates the success rate of a robot manipulating different objects when the scenario is fixed.  
235 This test set contains all the objects introduced during the data generation process of the  $i$ -th iteration.  
236 A higher value of this metric not only indicates better object generalization performance of the policy  
237 but also implies a better capability to enhance the diversity of objects in data generation. All success  
238 rates reported as mean values from 5 independent runs. Detailed compositions of all evaluation sets  
239 and success rate calculation method are provided in Appendix A.4.

240 **Baselines.** We compared our approach against the following methods: (1) Human Demo (Default):  
241 20 human demonstrations per task in a fixed scenario; (2) Human Demo (Enhanced): 20 demonstra-  
242 tions collected across diverse environments; (3) DexMimicGen(Default) [10]: generates data from a  
243 single demonstration via trajectory replay and editing; and (4) DexMimicGen (Enhanced): To create  
244 the strongest possible baseline, we provided DexMimicGen with 10 diverse human demonstrations  
245 collected across different scenarios—a 10x resource advantage compared to our method’s starting  
246 point. (4) w/o Res: An ablation model maintaining the base policy and residual policy while removing  
247 the  $\mathcal{A}_{EP}$  module. (5) w/o  $\mathcal{A}_{EP}$ : An ablation model featuring the base policy with  $\mathcal{A}_{EP}$  but excluding  
248 the Residual Reinforcement Learning component. (6) w/o Res. + w/o  $\mathcal{A}_{EP}$ : The minimal baseline  
249 configuration consisting solely of the base network, without either residual policy or  $\mathcal{A}_{EP}$  components.  
250 All methods are evaluated under identical conditions: same source demonstrations, simulation setups,  
251 model architectures, and test sets, ensuring a fair and rigorous comparison.

## 252 5.2 Data Diversity and Policy Generalization Performance Evaluation

253 As shown in the mid columns of Table 1, DexFlyWheel successfully expands data diversity with  
254 each iteration. In the final iteration ( $i=3$ ), our method generates an average of 2,250 various sce-  
255 nario configurations spanning 21 different objects per task—*all starting from just a single human*  
256 *demonstration per task*. These results underscore our method’s exceptional capability for automated  
257 generation of large-scale, diverse, and high-quality demonstration data, significantly reducing the  
258 dependency on extensive human data collection. The object diversity results (shown in the  $O$  column  
259 of Table 1) demonstrate our framework’s capacity for object-level data generation. Across all tasks,  
260 we expand from a single object per task in iteration 1 to an average of 20 distinct objects per task  
261 in iteration 3. This improvement stems from our residual reinforcement learning approach, which  
262 enhances adaptation to novel objects during the data generation process. As shown in the  $SR$  in  $T_O^i$   
263 column, our residual policies consistently improve performance on object generalization by 28.8% on  
264 average (comparing base policies to combined policies with residual components). For multi-factor  
265 generalization performance (presented in the  $SR$  in  $T_{OEP}$  column of Table 1), we observe that  
266 as dataset diversity increases across iterations, the generalization capabilities of trained policies  
267 correspondingly improve, achieving an average success rate of 81.9% in iteration 3—a substantial  
268 improvement from the initial 16.5% in iteration 1.

## 269 5.3 Performance Comparison Between DexFlyWheel and Baselines

270 As shown in Table 2, DexFlyWheel significantly outperforms baseline methods across all four tasks.  
271 Our findings can be summarized as follows: (a) Our method consistently outperforms all baselines,  
272 with an average 28.7% absolute improvement over DexMimicGen (Enhanced), the strongest baseline  
273 approach. This substantial performance gap can be attributed to our systematic enhancement of  
274 demonstration diversity through the iterative data FlyWheel mechanism. (b) Replay-based data  
275 generation methods demonstrate high sensitivity to task complexity, often failing in sophisticated  
276 scenarios involving multi-arm coordination such as lift and handover tasks. This limitation likely  
277 stems from their offline replay nature, which operates without awareness of changing environmental  
278 dynamics and robot behavioral patterns. In contrast, DexFlyWheel maintains consistently high  
279 success rates even in these challenging tasks, highlighting the critical importance of its interactive  
280 data generation process in enhancing generalization and robustness.(detailed analysis in Section 5.4.)  
281 (b) Collecting 20 trajectories in identical scenarios (Human Demo) enables basic skill acquisition  
282 but results in poor generalization performance (below 17% on most tasks). While introducing  
283 limited diversity can moderately improve performance on simpler tasks, it proves insufficient or  
284 even counterproductive for complex manipulation tasks involving dual-arm coordination, such as  
285 handover and lift, where Human Demo (Enhanced) achieves merely 2.5% or 0% success rates. This

Table 1: **Self-improving Data Generation Process.** *Left)* Evaluated tasks and iterations. *Middle)* Diversity metrics of our generated datasets. *Right)* Generalization performance of policies trained on our datasets.

Settings		Data Diversity					Generalization Performance	
Task	Iter.	O	E	P	Configs	Traj.	SR in $T_O^i$	SR in $T_{OEP}$
<b>Grasp</b>	$i = 1$	1	3	5	15	20	100.0% $\pm$ 1.2%	15.0% $\pm$ 2.1%
	$i = 2$	11	8	10	880	100	71.0% $\pm$ 4.3% $\rightarrow$ <b>84.0%</b> $\pm$ 3.5%	58.0% $\pm$ 4.8%
	$i = 3$	<b>22</b>	<b>12</b>	<b>15</b>	<b>3960</b>	<b>500</b>	35.0% $\pm$ 5.2% $\rightarrow$ <b>89.1%</b> $\pm$ 3.9%	<b>90.0%</b> $\pm$ 3.2%
<b>Pour</b>	$i = 1$	1	7	3	21	20	100.0% $\pm$ 3.1%	36.1% $\pm$ 3.3%
	$i = 2$	4	9	8	288	100	58.3% $\pm$ 5.1% $\rightarrow$ <b>75.0%</b> $\pm$ 4.0%	55.6% $\pm$ 4.5%
	$i = 3$	<b>12</b>	<b>12</b>	<b>10</b>	<b>1440</b>	<b>500</b>	58.0% $\pm$ 4.8% $\rightarrow$ <b>80.7%</b> $\pm$ 3.7%	<b>85.8%</b> $\pm$ 3.5%
<b>Lift</b>	$i = 1$	1	1	1	1	20	90.0% $\pm$ 2.2%	13.9% $\pm$ 2.8%
	$i = 2$	6	5	2	60	100	50.0% $\pm$ 5.3% $\rightarrow$ <b>83.3%</b> $\pm$ 3.8%	44.4% $\pm$ 4.6%
	$i = 3$	<b>26</b>	<b>12</b>	<b>5</b>	<b>1560</b>	<b>500</b>	68.8% $\pm$ 4.4% $\rightarrow$ <b>98.0%</b> $\pm$ 2.1%	<b>79.4%</b> $\pm$ 7.9%
<b>Handover</b>	$i = 1$	1	1	1	1	20	28.8% $\pm$ 4.7%	0.8% $\pm$ 1.1%
	$i = 2$	6	5	2	60	100	28.6% $\pm$ 5.8% $\rightarrow$ <b>85.7%</b> $\pm$ 4.2%	17.5% $\pm$ 3.4%
	$i = 3$	<b>20</b>	<b>12</b>	<b>5</b>	<b>1200</b>	<b>500</b>	32.1% $\pm$ 5.5% $\rightarrow$ <b>62.5%</b> $\pm$ 4.3%	<b>72.5%</b> $\pm$ 4.1%
Avg. $i = 1$		1.0	3.0	2.5	9.5	20	79.7%	16.5%
Avg. $i = 2$		6.8	6.8	5.5	322.0	100	52.0% $\rightarrow$ <b>82.0%</b>	43.9%
Avg. $i = 3$		<b>20.0</b>	<b>12.0</b>	<b>8.8</b>	<b>2040.0</b>	<b>500</b>	48.5% $\rightarrow$ <b>82.6%</b>	<b>81.9%</b>
Improvement ( $i = 1 \rightarrow 3$ )		<b>20.0x</b>	<b>4.0x</b>	<b>3.5x</b>	<b>214.7x</b>	<b>25.0x</b>	—	+396.4%

Notes:  $O$ : Number of objects,  $E$ : Number of environments,  $P$ : Number of poses. Configs: Total scenarios configurations ( $O \times E \times P$ ). Rightward arrows ( $\rightarrow$ ) show improvement from base policy to combined policy via residual policy.

Table 2: **Comparison with Baselines.** Success rates of policies trained on datasets generated by different methods when tested on multi-factor generalization test set ( $T_{OEP}$ ).

Method	Grasp	Pour	Lift	Handover	Avg.
Human Demo (Default)	6.1% $\pm$ 1.2%	16.7% $\pm$ 2.5%	13.9% $\pm$ 2.1%	0.8% $\pm$ 1.1%	9.4%
Human Demo (Enhanced)	15.0% $\pm$ 2.1%	36.1% $\pm$ 3.3%	2.5% $\pm$ 1.1%	0% $\pm$ 0.0%	13.4%
DexMimicGen (Default)	30.3% $\pm$ 3.8%	38.9% $\pm$ 4.2%	28.2% $\pm$ 3.5%	28.3% $\pm$ 4.7%	31.4%
DexMimicGen (Enhanced)	50.3% $\pm$ 4.5%	44.4% $\pm$ 3.8%	43.7% $\pm$ 3.6%	42.5% $\pm$ 4.9%	45.2%
Ours	<b>90.0%</b> $\pm$ 3.2%	<b>85.8%</b> $\pm$ 3.5%	<b>79.4%</b> $\pm$ 7.9%	<b>72.5%</b> $\pm$ 4.1%	<b>81.9%</b>

underscores the necessity of our approach which balances diversity expansion with data quality maintenance.

#### 5.4 Performance Analysis

In this section, we conduct an in-depth analysis to uncover the underlying reasons behind DexFly-Wheel’s superior performance.

To illustrate our framework’s object generalization advantages, we analyze the number of successfully handled objects in the data generation process across different tasks for our method, Ours w/o Res, and DexMimicGen in Figure 4. This figure demonstrates that DexFlyWheel achieves substantially greater object diversity compared to DexMimicGen. This superior performance can be attributed primarily to the residual reinforcement learning module, as evidenced by the significant drop in performance when this component is removed (w/o Res vs. ours: from 8.25 to 20 objects on average). Replay-based methods like DexMimicGen are geometry-unaware, typically operating effectively only on geometrically similar objects (same categories and shapes), which severely limits their ability to handle diverse object sets. To further investigate the robustness of different data generation approaches, we evaluated the data collection success rate across all tasks. In Table 3, we find that DexMimicGen’s success rate decreases as task complexity increases, achieving only 14.8% success rate when generating handover demonstrations, whereas our method still maintains 85.7%. The underlying reason is that objects experience varying dynamics during manipulating. Despite same

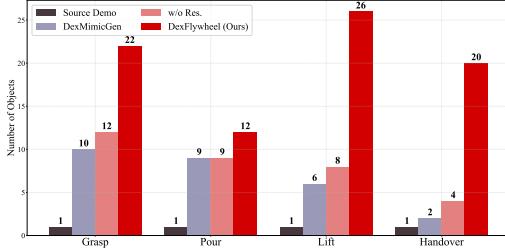


Figure 4: **Comparison of Object Diversity.** Our method successfully handles objects with diverse geometries, sizes, and categories.

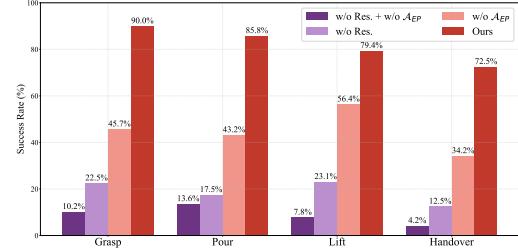


Figure 5: **Ablation Study.** Quantitative contribution of each module in DexFlyWheel across four manipulation tasks.

Table 3: **Data Collection Success Rate Comparison.** Success rates of generating valid demonstrations using different methods across tasks.

Method	Grasp	Pour	Lift	Handover
Human Demo	$100.0\% \pm 0.0\%$	$100.0\% \pm 0.0\%$	$100.0\% \pm 0.0\%$	$100.0\% \pm 0.0\%$
DexMimicGen	$87.3\% \pm 3.1\%$	$81.5\% \pm 4.3\%$	$68.2\% \pm 5.6\%$	$14.8\% \pm 6.2\%$
DexFlyWheel (Ours)	$93.6\% \pm 2.2\%$	$90.2\% \pm 2.8\%$	$89.5\% \pm 3.5\%$	$85.7\% \pm 5.1\%$

304 initial states, inherent task dynamics always exist, yet offline replay trajectories fail due to their  
305 dynamics-unaware characitcs.

306 We ablate the core components of DexFlyWheel across four manipulation tasks in Figure 5 We find  
307 that the residual policy contributes most significantly to generalization performance.

### 308 5.5 Deployment on Real Robots

309 We transfer our trained policy in simulation into real-world scenarios following a real-to-sim-to-real  
310 pipeline. Specifically, we collect one real-world demonstration via teleoperation and replay it in sim-  
311 ulation as the initial seed data. Then we generate hundreds of augmented successful demonstra-  
312 tions leveraging DexFlyWheel and distill them into a state-based diffusion policy. The diffusion policy  
313 takes the robot proprioception information and the object pose as input and output the robot joint  
314 action. Finally, we directly transfer our trained policy back to the real-world scenarios. We set up  
315 identical hardware settings both in simulation and real-world as shown in Figure 3, which includes  
316 two Realman RM75-6F arms paired with two PsiBot G0-R hands. We employ an egocentric-view  
317 RealSense D455 camera to obtain the real-world object pose leveraging FoundationPose [50]. We  
318 evaluate the performance of our pipeline on the Dual-arm Lift and Handover Task. The experimental  
319 results demonstrate that our simulation-trained policy can be successfully applied to real-world  
320 scenarios, achieving average success rates of 78.3% for the Dual-arm Lift task and 63.3% for the  
321 Handover task across three trials (20 trajectories per trial).

## 322 6 Limitations and Future Work

323 There are several limitations of our work. First, our reinforcement learning process still requires  
324 manually designed rewards, which can be time-consuming and may introduce biases. Future research  
325 could explore leveraging large language models to automate reward generation for diverse manipula-  
326 tion tasks. Second, our policies and simulations currently lack tactile feedback due to the immaturity  
327 of tactile sensing and simulation technologies. We plan to explore the potential of sensor-based tactile  
328 signals for contact-rich tasks in future work.

## 329 7 Conclusion

330 We present DexFlyWheel, a self-improving framework for generating diverse, high-quality dexterous  
331 manipulation data from minimal seed demonstrations. Our two-stage pipeline combines imitation  
332 learning’s behavioral priors with reinforcement learning’s generalization capabilities to expand data  
333 distribution across diverse environments. Experiments show policies trained on our generated data  
334 achieve 75% success rates in challenging scenarios, outperforming baselines, with successful transfer  
335 to real-world hardware.

336 **References**

- 337 [1] Stefan Schaal. Learning from demonstration. *Advances in neural information processing*  
338 *systems*, 9, 1996.
- 339 [2] Yunfei Bai and C. Karen Liu. Dexterous manipulation using both palm and fingers. In  
340 *International Conference on Robotics and Automation (ICRA)*, pages 1560–1565, 2014.
- 341 [3] S. Gruber. Robot hands and the mechanics of manipulation. *IEEE Journal on Robotics and*  
342 *Automation*, 2(1):59–59, 1986.
- 343 [4] Vikash Kumar, Yuval Tassa, Tom Erez, and Emanuel Todorov. Real-time behaviour synthesis for  
344 dynamic hand-manipulation. In *International Conference on Robotics and Automation (ICRA)*,  
345 pages 6808–6815, 2014.
- 346 [5] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo  
347 Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming  
348 Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang  
349 Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky.  $\pi_0$ : A  
350 vision-language-action flow model for general robot control. *CoRR*, abs/2410.24164, 2024.
- 351 [6] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski,  
352 Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu,  
353 Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander  
354 Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov,  
355 Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk  
356 Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo,  
357 Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspia Singh, Anikait Singh, Radu Soricuț,  
358 Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart,  
359 Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2:  
360 Vision-language-action models transfer web knowledge to robotic control, 2023.
- 361 [7] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu,  
362 Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *arXiv*  
363 *preprint arXiv:2410.07864*, 2024.
- 364 [8] Minjie Zhu, Yichen Zhu, Jinming Li, Junjie Wen, Zhiyuan Xu, Ning Liu, Ran Cheng, Chaomin  
365 Shen, Yixin Peng, Feifei Feng, and Jian Tang. Scaling diffusion policy in transformer to 1  
366 billion parameters for robotic manipulation. *arXiv preprint arXiv:2409.14411*, 2024.
- 367 [9] Chen Wang, Haochen Shi, Weizhuo Wang, Ruohan Zhang, Li Fei-Fei, and C. Karen Liu.  
368 Dexcap: Scalable and portable mocap data collection system for dexterous manipulation, 2024.
- 369 [10] Zhenyu Jiang, Yuqi Xie, Kevin Lin, Zhenjia Xu, Weikang Wan, Ajay Mandlekar, Linxi Fan,  
370 and Yuke Zhu. Dexmimicgen: Automated data generation for bimanual dexterous manipulation  
371 via imitation learning, 2024.
- 372 [11] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. Rlbench: The robot  
373 learning benchmark learning environment, 2019.
- 374 [12] Ajay Mandlekar, Soroush Nasiriany, Bowen Wen, Iretiayo Akinola, Yashraj Narang, Linxi Fan,  
375 Yuke Zhu, and Dieter Fox. Mimicgen: A data generation system for scalable robot learning  
376 using human demonstrations, 2023.
- 377 [13] Soroush Nasiriany, Abhiram Maddukuri, Lance Zhang, Adeet Parikh, Aaron Lo, Abhishek  
378 Joshi, Ajay Mandlekar, and Yuke Zhu. Robocasa: Large-scale simulation of everyday tasks for  
379 generalist robots, 2024.
- 380 [14] Yufei Wang, Zhou Xian, Feng Chen, Tsun-Hsuan Wang, Yian Wang, Katerina Fragkiadaki,  
381 Zackory Erickson, David Held, and Chuang Gan. Robogen: Towards unleashing infinite data  
382 for automated robot learning via generative simulation, 2024.

- 383 [15] Yao Mu, Tianxing Chen, Shijia Peng, Zanxin Chen, Zeyu Gao, Yude Zou, Lunkai Lin, Zhiqiang  
 384 Xie, and Ping Luo. Robotwin: Dual-arm robot benchmark with generative digital twins (early  
 385 version), 2024.
- 386 [16] Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: A large-  
 387 scale benchmark for general object grasping. In *Proceedings of the IEEE/CVF Conference on*  
 388 *Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- 389 [17] Fengshuo Bai, Yu Li, Jie Chu, Tawei Chou, Runchuan Zhu, Ying Wen, Yaodong Yang, and  
 390 Yuanpei Chen. Retrieval dexterity: Efficient object retrieval in clutters with dexterous hand.  
 391 *arXiv preprint arXiv:2502.18423*, 2025.
- 392 [18] Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang,  
 393 Stone Tao, Xinyue Wei, Yunchao Yao, et al. Maniskill2: A unified benchmark for generalizable  
 394 manipulation skills. *arXiv preprint arXiv:2302.04659*, 2023.
- 395 [19] Yuanpei Chen, Yiran Geng, Fangwei Zhong, Jiaming Ji, Jiechuang Jiang, Zongqing Lu, Hao  
 396 Dong, and Yaodong Yang. Bi-dexhands: Towards human-level bimanual dexterous manipulation.  
 397 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(5):2804–2818, 2023.
- 398 [20] Yuanpei Chen, Chen Wang, Li Fei-Fei, and C Karen Liu. Sequential dexterity: Chaining  
 399 dexterous policies for long-horizon manipulation. *arXiv preprint arXiv:2309.00987*, 2023.
- 400 [21] Ajay Mandlekar, Soroush Nasiriany, Bowen Wen, Iretiayo Akinola, Yashraj Narang, Linxi Fan,  
 401 Yuke Zhu, and Dieter Fox. Mimicgen: A data generation system for scalable robot learning  
 402 using human demonstrations. *arXiv preprint arXiv:2310.17596*, 2023.
- 403 [22] Matthew T Mason and J Kenneth Salisbury Jr. Robot hands and the mechanics of manipulation.  
 404 1985.
- 405 [23] Antonio Bicchi. Hands for dexterous manipulation and robust grasping: A difficult road toward  
 406 simplicity. *IEEE Transactions on robotics and automation*, 16(6):652–662, 2002.
- 407 [24] Igor Mordatch, Zoran Popović, and Emanuel Todorov. Contact-invariant optimization for hand  
 408 manipulation. In *Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer*  
 409 *animation*, pages 137–144, 2012.
- 410 [25] Vikash Kumar, Yuval Tassa, Tom Erez, and Emanuel Todorov. Real-time behaviour synthesis  
 411 for dynamic hand-manipulation. In *2014 IEEE International Conference on Robotics and*  
 412 *Automation (ICRA)*, pages 6808–6815. IEEE, 2014.
- 413 [26] Yunfei Bai and C Karen Liu. Dexterous manipulation using both palm and fingers. In *2014*  
 414 *IEEE International Conference on Robotics and Automation (ICRA)*, pages 1560–1565. IEEE,  
 415 2014.
- 416 [27] Yuanpei Chen, Tianhao Wu, Shengjie Wang, Xidong Feng, Jiechuang Jiang, Stephen Marcus  
 417 McAleer, Yiran Geng, Hao Dong, Zongqing Lu, Song-Chun Zhu, and Yaodong Yang. Towards  
 418 human-level bimanual dexterous manipulation with reinforcement learning, 2022.
- 419 [28] Chen Tang, Ben AbbateMatteo, Jiaheng Hu, Rohan Chandra, Roberto Martín-Martín, and Peter  
 420 Stone. Deep reinforcement learning for robotics: A survey of real-world successes, 2024.
- 421 [29] Shikhar Bahl, Abhinav Gupta, and Deepak Pathak. Human-to-robot imitation in the wild. *arXiv*  
 422 *preprint arXiv:2207.09450*, 2022.
- 423 [30] Priyanka Mandikal and Kristen Grauman. Dexvip: Learning dexterous grasping with human  
 424 hand pose priors from video. In *Conference on Robot Learning*, pages 651–661. PMLR, 2022.
- 425 [31] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances  
 426 from human videos as a versatile representation for robotics. In *Proceedings of the IEEE/CVF*  
 427 *Conference on Computer Vision and Pattern Recognition*, pages 13778–13790, 2023.
- 428 [32] Kenneth Shaw, Shikhar Bahl, and Deepak Pathak. Videodex: Learning dexterity from internet  
 429 videos. In *Conference on Robot Learning*, pages 654–665. PMLR, 2023.

- 430 [33] Sumedh Sontakke, Jesse Zhang, Séb Arnold, Karl Pertsch, Erdem Büyük, Dorsa Sadigh, Chelsea  
 431 Finn, and Laurent Itti. Roboclip: One demonstration is enough to learn robot policies. *Advances*  
 432 in Neural Information Processing Systems
- 433 [34] Robert McCarthy, Daniel CH Tan, Dominik Schmidt, Fernando Acero, Nathan Herr, Yilun Du,  
 434 Thomas G Thuruthel, and Zhibin Li. Towards generalist robot learning from internet video: A  
 435 survey. *arXiv preprint arXiv:2404.19664*, 2024.
- 436 [35] Arpit Bahety, Priyanka Mandikal, Ben Abbatematteo, and Roberto Martín-Martín. Screwmimic:  
 437 Bimanual imitation from human videos with screw space projection. *arXiv preprint*  
 438 *arXiv:2405.03666*, 2024.
- 439 [36] Hanzhi Chen, Boyang Sun, Anran Zhang, Marc Pollefeys, and Stefan Leutenegger. Vidbot:  
 440 Learning generalizable 3d actions from in-the-wild 2d human videos for zero-shot robotic  
 441 manipulation. *arXiv preprint arXiv:2503.07135*, 2025.
- 442 [37] Lai Sum Yim, Quang TN Vo, Ching-I Huang, Chi-Ruei Wang, Wren McQueary, Hsueh-Cheng  
 443 Wang, Haikun Huang, and Lap-Fai Yu. Wfh-vr: Teleoperating a robot arm to set a dining table  
 444 across the globe via virtual reality. In *2022 IEEE/RSJ International Conference on Intelligent*  
 445 *Robots and Systems (IROS)*, pages 4927–4934. IEEE, 2022.
- 446 [38] Yuzhe Qin, Wei Yang, Binghao Huang, Karl Van Wyk, Hao Su, Xiaolong Wang, Yu-Wei Chao,  
 447 and Dieter Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation  
 448 system. *arXiv preprint arXiv:2307.04577*, 2023.
- 449 [39] Xuxin Cheng, Jialong Li, Shiqi Yang, Ge Yang, and Xiaolong Wang. Open-television: Teleop-  
 450 eration with immersive active visual feedback. *arXiv preprint arXiv:2407.01512*, 2024.
- 451 [40] Runyu Ding, Yuzhe Qin, Jiyue Zhu, Chengzhe Jia, Shiqi Yang, Ruihan Yang, Xiaojuan Qi, and  
 452 Xiaolong Wang. Bunny-visionpro: Real-time bimanual dexterous teleoperation for imitation  
 453 learning. *arXiv preprint arXiv:2407.03162*, 2024.
- 454 [41] Jialiang Zhang, Haoran Liu, Danshi Li, Xinqiang Yu, Haoran Geng, Yufei Ding, Jiayi Chen,  
 455 and He Wang. Dexgraspnet 2.0: Learning generative dexterous grasping in large-scale synthetic  
 456 cluttered scenes, 2024.
- 457 [42] Pu Hua, Minghuan Liu, Annabella Macaluso, Yunfeng Lin, Weinan Zhang, Huazhe Xu, and  
 458 Lirui Wang. Gensim2: Scaling robot data generation with multi-modal and reasoning llms,  
 459 2024.
- 460 [43] Huy Ha, Pete Florence, and Shuran Song. Scaling up and distilling down: Language-guided  
 461 robot skill acquisition. In *Conference on Robot Learning*, pages 3766–3777. PMLR, 2023.
- 462 [44] Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan Zhang, and Li Fei-Fei. Rekep: Spatio-  
 463 temporal reasoning of relational keypoint constraints for robotic manipulation. *arXiv preprint*  
 464 *arXiv:2409.01652*, 2024.
- 465 [45] Shyam Sundar Kannan, Vishnunandan L. N. Venkatesh, and Byung-Cheol Min. Smart-lm:  
 466 Smart multi-agent robot task planning using large language models, 2024.
- 467 [46] Yunfei Li, Ying Yuan, Jingzhi Cui, Haoran Huan, Wei Fu, Jiaxuan Gao, Zekai Xu, and Yi Wu.  
 468 Robot generating data for learning generalizable visual robotic manipulation. In *2024 IEEE/RSJ*  
 469 *International Conference on Intelligent Robots and Systems (IROS)*, pages 5813–5820. IEEE,  
 470 2024.
- 471 [47] Zhengrong Xue, Shuying Deng, Zhenyang Chen, Yixuan Wang, Zhecheng Yuan, and Huazhe  
 472 Xu. Demogen: Synthetic demonstration generation for data-efficient visuomotor policy learning.  
 473 *arXiv preprint arXiv:2502.16932*, 2025.
- 474 [48] Xiu Yuan, Tongzhou Mu, Stone Tao, Yunhao Fang, Mengke Zhang, and Hao Su. Policy decor-  
 475 ator: Model-agnostic online refinement for large policy model. *arXiv preprint arXiv:2412.13630*,  
 476 2024.

- 477 [49] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-  
478 Martín, Chen Wang, Gabrael Levine, Michael Lingelbach, Jiankai Sun, et al. Behavior-1k:  
479 A benchmark for embodied ai with 1,000 everyday activities and realistic simulation. In  
480 Conference on Robot Learning, pages 80–93. PMLR, 2023.
- 481 [50] Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield. Foundationpose: Unified 6d pose  
482 estimation and tracking of novel objects. [arXiv preprint arXiv:2312.08344](#), 2023.

483 **NeurIPS Paper Checklist**

484 **1. Claims**

485 Question: Do the main claims made in the abstract and introduction accurately reflect the  
486 paper's contributions and scope?

487 Answer: [Yes]

488 Justification: The claims in the abstract and introduction match the core contributions and  
489 are supported by the method and results sections.

490 Guidelines:

- 491 • The answer NA means that the abstract and introduction do not include the claims  
492 made in the paper.
- 493 • The abstract and/or introduction should clearly state the claims made, including the  
494 contributions made in the paper and important assumptions and limitations. A No or  
495 NA answer to this question will not be perceived well by the reviewers.
- 496 • The claims made should match theoretical and experimental results, and reflect how  
497 much the results can be expected to generalize to other settings.
- 498 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
499 are not attained by the paper.

500 **2. Limitations**

501 Question: Does the paper discuss the limitations of the work performed by the authors?

502 Answer: [Yes]

503 Justification: The paper includes a dedicated discussion of limitations.

504 Guidelines:

- 505 • The answer NA means that the paper has no limitation while the answer No means that  
506 the paper has limitations, but those are not discussed in the paper.
- 507 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 508 • The paper should point out any strong assumptions and how robust the results are to  
509 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
510 model well-specification, asymptotic approximations only holding locally). The authors  
511 should reflect on how these assumptions might be violated in practice and what the  
512 implications would be.
- 513 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
514 only tested on a few datasets or with a few runs. In general, empirical results often  
515 depend on implicit assumptions, which should be articulated.
- 516 • The authors should reflect on the factors that influence the performance of the approach.  
517 For example, a facial recognition algorithm may perform poorly when image resolution  
518 is low or images are taken in low lighting. Or a speech-to-text system might not be  
519 used reliably to provide closed captions for online lectures because it fails to handle  
520 technical jargon.
- 521 • The authors should discuss the computational efficiency of the proposed algorithms  
522 and how they scale with dataset size.
- 523 • If applicable, the authors should discuss possible limitations of their approach to  
524 address problems of privacy and fairness.
- 525 • While the authors might fear that complete honesty about limitations might be used by  
526 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
527 limitations that aren't acknowledged in the paper. The authors should use their best  
528 judgment and recognize that individual actions in favor of transparency play an impor-  
529 tant role in developing norms that preserve the integrity of the community. Reviewers  
530 will be specifically instructed to not penalize honesty concerning limitations.

531 **3. Theory assumptions and proofs**

532 Question: For each theoretical result, does the paper provide the full set of assumptions and  
533 a complete (and correct) proof?

534 Answer: [NA]

535 Justification: The paper does not include formal theoretical results or proofs, as it focuses  
536 on algorithm design and empirical evaluation.

537 Guidelines:

- 538 • The answer NA means that the paper does not include theoretical results.
- 539 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
540 referenced.
- 541 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 542 • The proofs can either appear in the main paper or the supplemental material, but if  
543 they appear in the supplemental material, the authors are encouraged to provide a short  
544 proof sketch to provide intuition.
- 545 • Inversely, any informal proof provided in the core of the paper should be complemented  
546 by formal proofs provided in appendix or supplemental material.
- 547 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 548 4. Experimental result reproducibility

549 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
550 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
551 of the paper (regardless of whether the code and data are provided or not)?

552 Answer: [Yes]

553 Justification: The paper provides detailed descriptions of the model architecture, training  
554 procedure, dataset, and evaluation metrics, enabling reproduction of the main results without  
555 requiring access to the code.

556 Guidelines:

- 557 • The answer NA means that the paper does not include experiments.
- 558 • If the paper includes experiments, a No answer to this question will not be perceived  
559 well by the reviewers: Making the paper reproducible is important, regardless of  
560 whether the code and data are provided or not.
- 561 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
562 to make their results reproducible or verifiable.
- 563 • Depending on the contribution, reproducibility can be accomplished in various ways.  
564 For example, if the contribution is a novel architecture, describing the architecture fully  
565 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
566 be necessary to either make it possible for others to replicate the model with the same  
567 dataset, or provide access to the model. In general, releasing code and data is often  
568 one good way to accomplish this, but reproducibility can also be provided via detailed  
569 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
570 of a large language model), releasing of a model checkpoint, or other means that are  
571 appropriate to the research performed.
- 572 • While NeurIPS does not require releasing code, the conference does require all submis-  
573 sions to provide some reasonable avenue for reproducibility, which may depend on the  
574 nature of the contribution. For example
  - 575 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
576 to reproduce that algorithm.
  - 577 (b) If the contribution is primarily a new model architecture, the paper should describe  
578 the architecture clearly and fully.
  - 579 (c) If the contribution is a new model (e.g., a large language model), then there should  
580 either be a way to access this model for reproducing the results or a way to reproduce  
581 the model (e.g., with an open-source dataset or instructions for how to construct  
582 the dataset).
  - 583 (d) We recognize that reproducibility may be tricky in some cases, in which case  
584 authors are welcome to describe the particular way they provide for reproducibility.  
585 In the case of closed-source models, it may be that access to the model is limited in  
586 some way (e.g., to registered users), but it should be possible for other researchers  
587 to have some path to reproducing or verifying the results.

#### 588 5. Open access to data and code

589 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
590 tions to faithfully reproduce the main experimental results, as described in supplemental  
591 material?

592 Answer: [Yes]

593 Justification: The code are provided in supplementary materials for reproduction of all key  
594 experiments, as detailed in the supplemental material.

595 Guidelines:

- 596 • The answer NA means that paper does not include experiments requiring code.
- 597 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 598 • While we encourage the release of code and data, we understand that this might not be  
600 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
601 including code, unless this is central to the contribution (e.g., for a new open-source  
602 benchmark).
- 603 • The instructions should contain the exact command and environment needed to run to  
604 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 605 • The authors should provide instructions on data access and preparation, including how  
606 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 607 • The authors should provide scripts to reproduce all experimental results for the new  
608 proposed method and baselines. If only a subset of experiments are reproducible, they  
609 should state which ones are omitted from the script and why.
- 610 • At submission time, to preserve anonymity, the authors should release anonymized  
611 versions (if applicable).
- 612 • Providing as much information as possible in supplemental material (appended to the  
613 paper) is recommended, but including URLs to data and code is permitted.

## 615 6. Experimental setting/details

616 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
617 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
618 results?

619 Answer: [Yes]

620 Justification: The paper specifies all relevant training details, including hyperparameters,  
621 optimizer settings, and training schedules, with additional configurations provided in the  
622 appendix.

623 Guidelines:

- 624 • The answer NA means that the paper does not include experiments.
- 625 • The experimental setting should be presented in the core of the paper to a level of detail  
626 that is necessary to appreciate the results and make sense of them.
- 627 • The full details can be provided either with the code, in appendix, or as supplemental  
628 material.

## 629 7. Experiment statistical significance

630 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
631 information about the statistical significance of the experiments?

632 Answer: [Yes]

633 Justification: We report the standard deviation over multiple runs to reflect variability, and  
634 clearly indicate it in both the main text and figure captions.

635 Guidelines:

- 636 • The answer NA means that the paper does not include experiments.
- 637 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
638 dence intervals, or statistical significance tests, at least for the experiments that support  
639 the main claims of the paper.

- 640           • The factors of variability that the error bars are capturing should be clearly stated (for  
 641           example, train/test split, initialization, random drawing of some parameter, or overall  
 642           run with given experimental conditions).  
 643           • The method for calculating the error bars should be explained (closed form formula,  
 644           call to a library function, bootstrap, etc.)  
 645           • The assumptions made should be given (e.g., Normally distributed errors).  
 646           • It should be clear whether the error bar is the standard deviation or the standard error  
 647           of the mean.  
 648           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
 649           preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
 650           of Normality of errors is not verified.  
 651           • For asymmetric distributions, the authors should be careful not to show in tables or  
 652           figures symmetric error bars that would yield results that are out of range (e.g. negative  
 653           error rates).  
 654           • If error bars are reported in tables or plots, The authors should explain in the text how  
 655           they were calculated and reference the corresponding figures or tables in the text.

656           **8. Experiments compute resources**

657           Question: For each experiment, does the paper provide sufficient information on the com-  
 658           puter resources (type of compute workers, memory, time of execution) needed to reproduce  
 659           the experiments?

660           Answer: [Yes]

661           Justification: The paper specifies the use of NVIDIA RTX4090 GPUs and 4 CPU cores,  
 662           along with details on training time and batch sizes for all experiments.

663           Guidelines:

- 664           • The answer NA means that the paper does not include experiments.  
 665           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
 666           or cloud provider, including relevant memory and storage.  
 667           • The paper should provide the amount of compute required for each of the individual  
 668           experimental runs as well as estimate the total compute.  
 669           • The paper should disclose whether the full research project required more compute  
 670           than the experiments reported in the paper (e.g., preliminary or failed experiments that  
 671           didn't make it into the paper).

672           **9. Code of ethics**

673           Question: Does the research conducted in the paper conform, in every respect, with the  
 674           NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

675           Answer: [Yes]

676           Justification: The research adheres to the NeurIPS Code of Ethics, with no ethical concerns  
 677           related to data usage, human subjects, or potential misuse identified.

678           Guidelines:

- 679           • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
 680           • If the authors answer No, they should explain the special circumstances that require a  
 681           deviation from the Code of Ethics.  
 682           • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
 683           eration due to laws or regulations in their jurisdiction).

684           **10. Broader impacts**

685           Question: Does the paper discuss both potential positive societal impacts and negative  
 686           societal impacts of the work performed?

687           Answer: [Yes]

688           Justification: The paper considers positive impacts in assistive and robotics, while also noting  
 689           possible misuse in surveillance and unintended physical harm from inaccurate predictions.

690           Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper does not release models or data with high risk for misuse; thus, no special safeguards are necessary.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: The paper does not use third-party assets that require attribution or license disclosure.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 744           • For scraped data from a particular source (e.g., website), the copyright and terms of  
745           service of that source should be provided.  
746           • If assets are released, the license, copyright information, and terms of use in the  
747           package should be provided. For popular datasets, [paperswithcode.com/datasets](http://paperswithcode.com/datasets)  
748           has curated licenses for some datasets. Their licensing guide can help determine the  
749           license of a dataset.  
750           • For existing datasets that are re-packaged, both the original license and the license of  
751           the derived asset (if it has changed) should be provided.  
752           • If this information is not available online, the authors are encouraged to reach out to  
753           the asset's creators.

754           **13. New assets**

755           Question: Are new assets introduced in the paper well documented and is the documentation  
756           provided alongside the assets?

757           Answer: [NA]

758           Justification: This paper does not release new assets.

759           Guidelines:

- 760           • The answer NA means that the paper does not release new assets.  
761           • Researchers should communicate the details of the dataset/code/model as part of their  
762           submissions via structured templates. This includes details about training, license,  
763           limitations, etc.  
764           • The paper should discuss whether and how consent was obtained from people whose  
765           asset is used.  
766           • At submission time, remember to anonymize your assets (if applicable). You can either  
767           create an anonymized URL or include an anonymized zip file.

768           **14. Crowdsourcing and research with human subjects**

769           Question: For crowdsourcing experiments and research with human subjects, does the paper  
770           include the full text of instructions given to participants and screenshots, if applicable, as  
771           well as details about compensation (if any)?

772           Answer: [NA]

773           Justification: This paper does not involve crowdsourcing.

774           Guidelines:

- 775           • The answer NA means that the paper does not involve crowdsourcing nor research with  
776           human subjects.  
777           • Including this information in the supplemental material is fine, but if the main contribu-  
778           tion of the paper involves human subjects, then as much detail as possible should be  
779           included in the main paper.  
780           • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
781           or other labor should be paid at least the minimum wage in the country of the data  
782           collector.

783           **15. Institutional review board (IRB) approvals or equivalent for research with human  
784           subjects**

785           Question: Does the paper describe potential risks incurred by study participants, whether  
786           such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
787           approvals (or an equivalent approval/review based on the requirements of your country or  
788           institution) were obtained?

789           Answer: [NA]

790           Justification: This paper does not involve crowdsourcing.

791           Guidelines:

- 792           • The answer NA means that the paper does not involve crowdsourcing nor research with  
793           human subjects.

- 794           • Depending on the country in which research is conducted, IRB approval (or equivalent)  
795           may be required for any human subjects research. If you obtained IRB approval, you  
796           should clearly state this in the paper.  
797           • We recognize that the procedures for this may vary significantly between institutions  
798           and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
799           guidelines for their institution.  
800           • For initial submissions, do not include any information that would break anonymity (if  
801           applicable), such as the institution conducting the review.

802           **16. Declaration of LLM usage**

803           Question: Does the paper describe the usage of LLMs if it is an important, original, or  
804           non-standard component of the core methods in this research? Note that if the LLM is used  
805           only for writing, editing, or formatting purposes and does not impact the core methodology,  
806           scientific rigorousness, or originality of the research, declaration is not required.

807           Answer: [NA]

808           Justification: The core method development in this research does not involve LLMs.

809           Guidelines:

- 810           • The answer NA means that the core method development in this research does not  
811           involve LLMs as any important, original, or non-standard components.  
812           • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
813           for what should or should not be described.