

Constraint-Preserving Data Generation for One-Shot Visuomotor Policy Generalization

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Large-scale demonstration data has powered key breakthroughs in
2 robot manipulation, but collecting that data remains costly and time-consuming.
3 To this end, we present Constraint-Preserving Data Generation (CP-Gen), a
4 method that uses a single expert trajectory to generate robot demonstrations con-
5 taining novel object geometries and poses. These generated demonstrations are
6 used to train closed-loop visuomotor policies that transfer zero-shot to the real
7 world. Similar to prior data-generation work focused on pose variations, CP-Gen
8 first decomposes expert demonstrations into free-space motions and robot skills.
9 Unlike prior work, we achieve geometry-aware data generation by formulating
10 robot skills as keypoint-trajectory constraints: keypoints on the robot or grasped
11 object must track a reference trajectory defined relative to a task-relevant object.
12 To generate a new demonstration, CP-Gen samples pose and geometry transforms
13 for each task-relevant object, then applies these transforms to the object and its as-
14 sociated keypoints or keypoint trajectories. We optimize robot joint configurations
15 so that the keypoints on the robot or grasped object track the transformed keypoint
16 trajectory, and then motion plan a collision-free path to the first optimized joint
17 configuration. Using demonstrations generated by CP-Gen, we train visuomotor
18 policies that generalize across variations in object geometries and poses. Experi-
19 ments on 16 simulation tasks and four real-world tasks, featuring multi-stage,
20 non-prehensile and tight-tolerance manipulation, show that policies trained using
21 our method achieve an average success rate of 77%, outperforming the best base-
22 line which achieves an average success rate of 50%. Simulation and real-world
23 videos can be found at: cpgen-anon.github.io.

24 **Keywords:** Imitation Learning, Data Generation, Robot Manipulation

25 1 Introduction

26 Teaching robots to act by imitating humans is one of the most powerful, yet costly, ideas in robotics.
27 At its best, large-scale imitation learning has enabled robots to tie shoelaces, fix broken grippers, and
28 even cook shrimp [1–4]. Further scaling promises to unlock broad generalization across complex
29 manipulation tasks. Behind these successes, however, lies a steep cost: human labor, months of
30 continuous robot operation, and heavy infrastructure demands. For example, ALOHA Unleashed [3]
31 collected 26,000 demonstrations over 8 months across 10 robots and 35 operators, while DROID [5]
32 gathered 76,000 real-world demonstrations over a year. As the field pushes forward, reducing the
33 burden of demonstration collection will be a key challenge.

34 Automated data generation promises to reduce the manual effort required in data collection. One
35 line of work [6–10] generates demonstrations for scenes with diverse object poses by collecting a
36 handful of teleoperated demonstrations and algorithmically multiplying them using pose transfor-
37 mations and action replay. Conceptually, these works aim to enable data efficient spatial generalization
38 by exploiting SE(3) equivariance of robot actions with respect to object poses [9]. However, be-
39 cause they rely on SE(3) transformations, these works cannot generate demonstrations that adapt
40 to geometric variations (such as different aspect ratios) with a given object instance. For example,

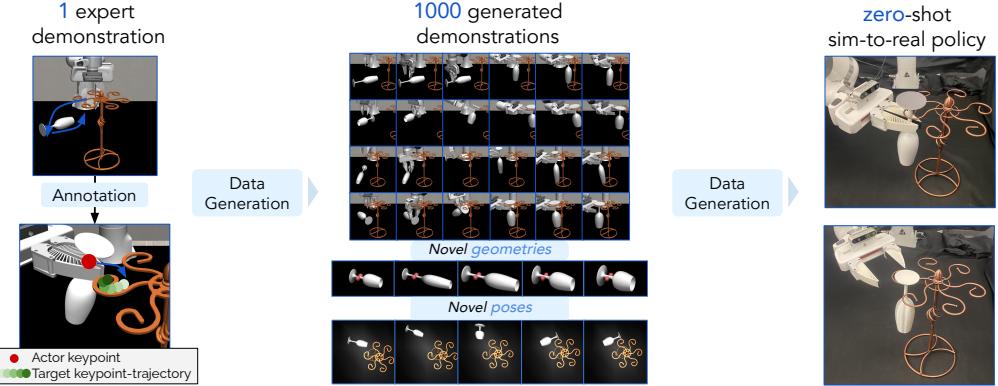


Figure 1: **CP-Gen** uses one expert demonstration and keypoint-trajectory constraints to generate diverse demonstrations involving novel object geometries and poses, enabling large-scale policy training and zero-shot sim-to-real transfer. For the Wine Glass Spiral Hanging task, the first keypoint-trajectory is anchored to the wine glass and constrains the motion of points on the end effector. The second keypoint-trajectory is anchored relative to the spiral-shaped wine glass rack, and constrains a keypoint on the wine glass stem to track a complex spiral motion.

41 robot actions that succeed in hanging a short, wide wine glass onto a rack may fail entirely for a tall,
42 narrow one, even if both are aligned under the same pose transformation.

43 A complementary strategy to ease the burden on manual data collection is to increase learning efficiency
44 by embedding structure—specifically, equivariance—into policy architectures [11–23]. One line of work focuses on open-loop manipulation, exploiting symmetry to generalize across different
45 object poses [24–27]. More recent work has moved toward closed-loop policy learning, embedding
46 equivariances directly into the policy architecture and directly predicting robot actions [11, 28–32].
47 However, embedding equivariance bottlenecks can limit policy expressiveness and be insufficient
48 when dealing with scenes that involve multiple interacting objects or symmetry-breaking variations.

50 In this work, we propose **Constraint-Preserving Data Generation (CP-Gen)**, a geometry-aware
51 data generation method that takes a single expert demonstration and generates new demonstrations
52 with both novel object poses and novel geometries. With CP-Gen, we aim to achieve the benefits of
53 equivariant policy architectures without the bottleneck of embedding structure directly into policies.

54 Our key insight is to bring geometry awareness into the data generation process through a **keypoint-
55 trajectory constraint** formulation. Starting from a *single* expert demonstration, we segment the
56 trajectory into *free-space motions* (segments that can be replaced with point-to-point collision-free
57 motion planning) and *robot skills* (segments requiring some interaction with task-relevant objects
58 that cannot be completed through pure collision-free motion planning). We then formulate each
59 skill as a keypoint-trajectory constraint: selected robot or grasped-object keypoints must track a
60 reference trajectory defined in the frame of a task-relevant object. For example, in the *Wine Glass*
61 *Spiral Hanging* task (Figure 1), the grasping skill is formulated as a trajectory-tracking constraint
62 between keypoints on the gripper tips and a target trajectory in the wine-glass frame, where the target
63 trajectory includes the pre-grasp segment up to the grasp pose. The subsequent *spiral* insertion
64 phase introduces a second constraint: keypoints on the glass stem must track a spiral-like trajectory
65 defined relative to the rack. By anchoring keypoints to the reference frame of task-relevant objects,
66 we can sample new object geometries and poses, adapt the keypoint-trajectory constraints according
67 to the transforms applied to the original object geometry, and generate thousands of demonstrations
68 covering diverse shapes and geometries. Using these demonstrations, we can train a visuomotor policy
69 that naturally generalizes robustly across changes in object pose and geometry, without requiring
70 additional human demonstrations.

71 Our contributions are threefold. First, we propose CP-Gen, a data generation method that produces
72 robot demonstrations from a single expert demonstration using a keypoint-trajectory constraint for-
73 mulation. Unlike prior work, and due to this formulation, CP-Gen generates data in a geometry
74 aware manner. Second, we validate our keypoint-trajectory data generation formulation on the Mim-
75 icGen simulation benchmark. We achieve state-of-the-art performance (85% vs. 63%) on the default

76 task settings, as well as on our proposed simulation benchmark featuring novel object geometries
77 (70% vs. 37%). Third, we demonstrate successful zero-shot sim-to-real on a set of four real-world
78 tasks featuring multi-stage, non-prehensile and tight-tolerance manipulation.

79 2 Related Work

80 **Data Generation for Robotics.** One approach to addressing data scarcity in robotics is to gen-
81 erate new demonstrations from teleoperated source data [6–8, 10, 33]. MimicGen [6] generates
82 demonstrations by resetting object poses and replaying pose-transformed versions of the source tra-
83 jectory. IntervenGen [33] extends MimicGen with human interventions mid-trajectory to correct
84 failures, while DexMimicGen [8] adapts MimicGen for bimanual and dexterous settings. Skill-
85 Gen [7] segments demonstrations into skills and motions, alternating between learned policies and
86 motion planning to improve robustness. At the core of these methods is a relative-pose formulation
87 and the exploitation of SE(3)-equivariance [9]: given an SE(3) transformation applied to an object,
88 they apply the same SE(3) transformation to robot actions to replicate the original effect. As a result,
89 they particularly struggle to generate demonstrations with new object geometries for tasks requiring
90 tight-tolerance manipulation. CP-Gen addresses this limitation by introducing a keypoint-trajectory
91 constraint formulation, representing skills in terms of robot or grasped-object keypoints tracking a
92 reference trajectory. This formulation enables demonstration generation not only across novel object
93 poses, but also across significantly different object geometries.

94 **Equivariance for Robotic Manipulation.** Another strategy to improve generalization is to em-
95 bed geometric structure into policy inference and architecture. Open-loop methods design SE(2)
96 or SE(3) equivariant feature representations on point clouds or images [11–23], achieving strong
97 sample efficiency but often relying on expensive inference-time optimization [24–27] and lack-
98 ing reactivity to dynamic changes. To address these limitations, recent work has moved toward
99 closed-loop policies, beginning with planar SO(2) symmetries and extending to full SIM(3) equiv-
100 ariance [11, 28–32], where SIM(3) captures rotation, translation, and uniform scaling transforma-
101 tions. While SIM(3)-equivariant policies enable more flexible closed-loop control compared to
102 SE(3)-equivariant policies, they still fundamentally assume that object geometry variations can be
103 captured by uniform scaling, and for larger geometry variations, they rely on policy generalization.
104 These closed-loop equivariance works [30–32] also tend to focus on single-object to robot equivari-
105 ance, instead of the multi-object scenario. In CP-Gen, rather than using architectural equivariance,
106 we leverage data generation to generalize to both novel poses and arbitrary geometry variations via
107 a keypoint-trajectory constraint formulation.

108 3 Problem Setting

109 We are given a single expert demonstration in the form of a trajectory $\tau_{\text{src}} = \{o_t, a_t\}_{t=1}^H$ consisting of
110 H observations o , and H actions a . Actions $a_t = \{a_{\text{eef}}, a_{\text{grip}}\}$ contain end-effector poses $a_{\text{eef}} \in \text{SE}(3)$
111 along with gripper actions a_{grip} for a given manipulation task. Observations $o_t = \{I_{\text{hand}}, I_{3\text{pv}}, q_t\}$ at
112 each time step consist of a wrist-mounted camera image I_{hand} , a third-person view camera image $I_{3\text{pv}}$,
113 and robot proprioception data q_t . We aim to train a visuomotor policy $\pi_\theta(a|o)$ mapping observations
114 $o \in O$ to actions $a \in A$.

115 Similar to prior data generation methods using source demonstrations [6, 7, 33], we assume that
116 an expert demonstration τ_{src} can be decomposed into a sequence of alternating free-space motion
117 segments $\tau_{\text{motion}}^{(i)}$ and robot skill segments $\tau_{\text{skill}}^{(i)}$: $\tau_{\text{src}} = [\tau_{\text{motion}}^{(1)}, \tau_{\text{skill}}^{(1)}, \tau_{\text{motion}}^{(2)}, \tau_{\text{skill}}^{(2)}, \dots]$. For each skill
118 segment, we additionally assume access to a set of annotated keypoints on the robot gripper or a
119 grasped object that are relevant to the completion of the skill. We refer to these as *actor keypoints*
120 and denote them by $\mathcal{K}_{\text{actor}}(t) = \{k_1^{\text{actor}}(t), \dots, k_N^{\text{actor}}(t)\}$, where each keypoint $k_i^{\text{actor}}(t) \in \mathbb{R}^3$ is defined
121 in the local frame of the robot gripper or the grasped object at time t . Although in this work we
122 manually selected these keypoints, in future work, keypoints can be predicted using vision-language
123 models [34–36], or inferred via task-specific keypoint detectors [37–40].

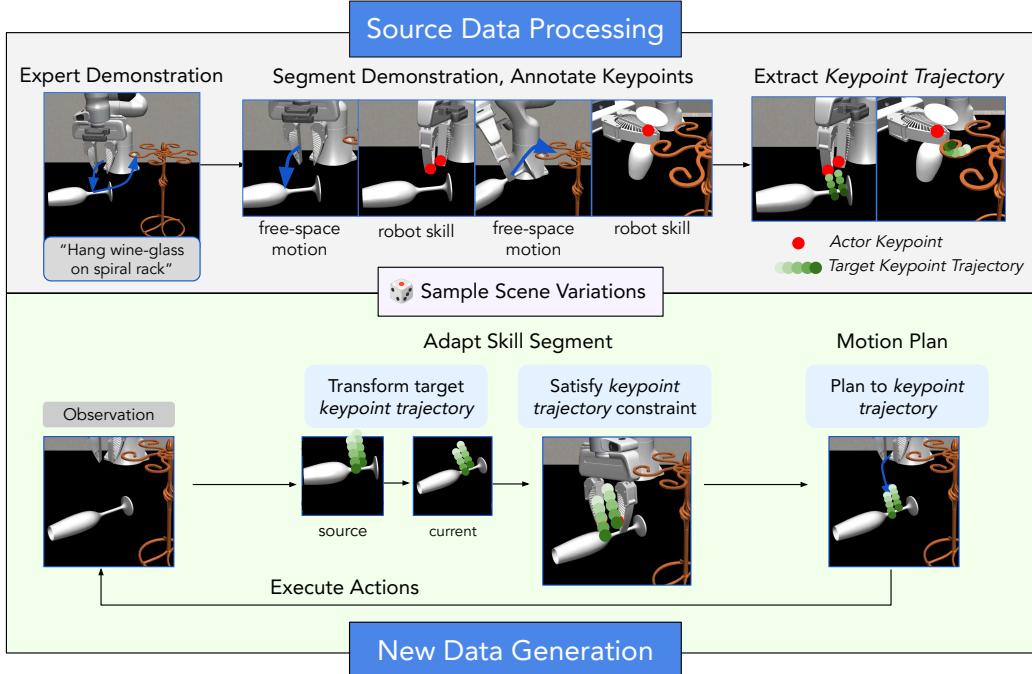


Figure 2: **CP-Gen Method.** In the *Source Data Processing* stage (top), starting from an expert demonstration τ_{src} , we (a) segment the trajectory into free-space motion and skill segments, (b) annotate keypoints on the robot or grasped object (dubbed *actor* keypoints), and (c) convert each skill segment into a *keypoint-trajectory constraint* by extracting a keypoint-trajectory expressed in the frame of a task-relevant object. After processing the source demonstration, in the *New Data Generation* stage (bottom), we (a) sample scene variations by applying geometry and pose transforms to every task-relevant object. For each sampled scene, we (a) adapt our source skill segment to the current observation by transforming the extracted target keypoint-trajectory with the current object transforms and solving for robot configurations that satisfy the updated keypoint-trajectory constraint (green dots), and (b) motion plan a collision-free path that from the current robot joint configuration to the first configuration in the solved configuration trajectory. The process iterates over all segments to generate a demonstration for the new scene, which can subsequently be used to train a visuomotor policy $\pi_\theta(a | o)$.

124 4 CP-Gen: Constraint-Preserving Data Generation

125 We present **Constraint-Preserving Data Generation (CP-Gen)**, a data generation method that
 126 enables one-shot visual imitation learning by adapting a source demonstration trajectory to scenes with
 127 novel object geometries and poses (see Figure 2). Our approach begins by decomposing the original
 128 demonstration trajectory τ_{src} into a sequence of skill segments τ_{skill} interleaved with collision-free
 129 motion segments τ_{motion} , such that $\tau_{src} = [\tau_{motion}, \tau_{skill}, \dots, \tau_{motion}, \tau_{skill}]$. Each skill segment τ_{skill}
 130 is formulated as a *keypoint-trajectory constraint*: keypoints on the robot or grasped object, *actor*
 131 *keypoints* ${}^O\mathcal{K}_{actor}(t)$, must track a reference trajectory defined relative to a task-relevant object. This
 132 reference trajectory is represented as a time-varying sequence of *target keypoints* ${}^O\mathcal{K}_{target}(t)$ ex-
 133 pressed in the local frame of the object O . We extract this reference keypoint-trajectory (Section 4.1)
 134 by transforming the actor keypoints into the object frame at each timestep. We assume that adapting
 135 and preserving this keypoint-trajectory constraint to a new scene leads to successful execution. To
 136 generate a new demonstration in a novel scene, we first sample *pose* and *geometry* transforms for
 137 each task-relevant object. We then apply the sampled transformations to the objects (and their asso-
 138 ciated keypoint trajectories) to produce a new scene reset. For the new scene, we iterate through two
 139 phases of skill segment adaptation (Section 4.2) and motion planning (Section 4.3). Finally, after
 140 filtering out any failed demonstrations using a success detector (which we assume exists for each
 141 task), we use imitation learning to train a visuomotor policy on the generated demonstrations.

142 4.1 Keypoint-Trajectory Constraint Extraction

143 To extract the target keypoint-trajectory for a skill segment τ_{skill} from the original demonstration
 144 τ_{src} , we transform the actor keypoints ${}^A\mathcal{K}_{actor} = \{{}^A k_1, \dots, {}^A k_N\}$, defined in the actor's local frame A

	StackThree	Square	Threading	Assembly	Kitchen	Coffee	MugCleanup	HammerCleanup
StackThreeG								
Property	Stack Three	Square	Threading	Assembly	Kitchen	Coffee	Mug Cleanup	Hammer Cleanup
Object(s)	block1, block2, block3	nut, peg	needle, tripod	base, piece_1, piece_2	bread, pot, stove, button	pod, machine	drawer, object	hammer, drawer
Scale Range	[0.6,1.4]	[1.0,1.0]	[0.9,1.1]	[0.7,1.3]	[0.7,1.4]	[0.6,1.1]	[0.7,1.3]	[0.7,1.2]

Figure 3: **Simulation Tasks and Geometry Generalization Variants.** *Top:* Tasks from the MimicGen benchmark [6] and object geometries sampled from our proposed *Geometry Generalization* task variants. *Bottom:* Uniform scale sampling ranges applied in *Geometry Generalization* task variants.

145 (robot gripper or grasped object), into the task-relevant object frame O via a frame transform from
146 actor to world frame, then world to task-relevant object frame:

$${}^O k_i(t) = \underbrace{{}^O T_W(t)}_{\text{object pose}^{-1}} \cdot \underbrace{{}^W T_A(t)}_{\text{actor pose}} \cdot {}^A k_i, \quad (1)$$

147 Applying Eq. 1 across all t yields target keypoint-trajectory: ${}^O \mathcal{K}_{\text{target}}(t) = \{{}^O k_1(t), \dots, {}^O k_N(t)\}$.

148 4.2 Skill Segment Adaptation

149 Given a skill segment τ_{skill} and a current scene observation, our goal is to adapt the original keypoint-
150 trajectory constraint to new task-relevant object geometries and poses. Doing so requires updating
151 both the actor keypoints and the target keypoint-trajectory.

152 **Update Keypoint-Trajectory Constraint for New Geometry.** Suppose that the task-relevant ob-
153 ject undergoes a geometric transformation \mathbf{X} (e.g., non-uniform scaling) in its own local frame O .
154 Let the original local-frame target keypoint-trajectory be ${}^O \mathcal{K}_{\text{target}}(t) = \{{}^O k_1(t), \dots, {}^O k_N(t)\}$. We
155 therefore apply \mathbf{X} to transform the target keypoint-trajectory:

$${}^O k'_i(t) = \mathbf{X} \cdot {}^O k_i(t), \quad \forall t.$$

156 If the actor keypoints ${}^A \mathcal{K}_{\text{actor}}(t) = \{{}^A k_1(t), \dots, {}^A k_N(t)\}$ are anchored to a grasped object that has
157 been geometrically transformed, we similarly apply the object’s local-frame geometric transformation
158 \mathbf{X}_A to the actor keypoints: ${}^A k'_i = \mathbf{X}_A \cdot {}^A k_i$. These updated keypoints ${}^A k'_i$ are then used in the
159 optimization process to match the newly transformed target trajectory.

160 **Solve for Robot Configuration via Keypoint Matching.** Given the updated actor keypoints ${}^A k'_i$ (in
161 the actor’s local frame) and the updated target keypoint-trajectory ${}^O k'_i(t)$ (in the object’s local frame),
162 our goal is to solve for the robot joint configuration q_t^* at each timestep t such that the transformed
163 actor keypoints match the target keypoints in the world frame.

164 First, we compute the world-frame position of each actor keypoint (recall that actor keypoints are
165 defined on the robot or a grasped object, and that we assume a fixed end-effector to grasped object
166 frame transformation) via forward kinematics: ${}^W k'_i(q) = f_{\text{FK}}(q, {}^A k'_i)$. Next, we compute the world-
167 frame position of each target keypoint by transforming the updated local-frame keypoints using the
168 current object pose ${}^W T_O(t)$ relative to the world frame: ${}^W k_i^{\text{target}}(t) = {}^W T_O(t) \cdot {}^O k'_i(t)$.

Method	StackThree	Square	Thread	Assembly	Kitchen	Coffee	Mug	Hammer	Avg
CP-Gen [RGB]	0.82±0.00	0.87±0.03	0.87±0.03	0.85±0.03	0.96±0.01	1.00±0.00	0.86±0.03	0.8±0.03	0.88
MimicGen [RGB]	0.98±0.01	0.72±0.04	0.38±0.04	0.32±0.04	0.98±0.01	0.58±0.04	0.78±0.03	0.64±0.04	0.67
CP-Gen [D+S]	0.56±0.00	0.91±0.02	0.85±0.03	0.79±0.03	0.92±0.02	0.99±0.01	0.72±0.00	0.79±0.03	0.82
MimicGen [D+S]	0.52±0.04	0.76±0.03	0.44±0.04	0.40±0.04	0.94±0.02	0.56±0.04	0.46±0.04	0.64±0.04	0.59

Method	StackThreeG	SquareG	ThreadG	AssemblyG	KitchenG	CoffeeG	MugG	HammerG	Avg
CP-Gen [RGB]	0.68±0.04	0.88±0.03	0.79±0.03	0.55±0.04	0.83±0.03	0.58±0.04	0.82±0.03	0.67±0.04	0.73
MimicGen [RGB]	0.60±0.04	0.42±0.04	0.10±0.02	0.16±0.03	0.74±0.04	0.00±0.00	0.38±0.04	0.40±0.04	0.35
CP-Gen [D+S]	0.36±0.00	0.83±0.03	0.73±0.04	0.79±0.03	0.86±0.03	0.63±0.04	0.76±0.04	0.74±0.04	0.67
MimicGen [D+S]	0.22±0.03	0.64±0.04	0.05±0.02	0.38±0.04	0.84±0.03	0.00±0.00	0.50±0.04	0.45±0.04	0.39

Table 1: **CP-Gen achieves state-of-the-art results on the MimicGen simulation benchmark.** On the original MimicGen benchmark (default task variants), CP-Gen achieves an average success rate of 88% compared to MimicGen’s 67%. On our custom benchmark containing *Geometry Generalization* task variants which feature novel object geometries (denoted as TaskG), CP-Gen achieves an average success rate of 70%, outperforming MimicGen’s 37% by a margin of 33%. These results highlight CP-Gen’s strong generalization not only to pose variations but also to challenging geometric variations. Bolded numbers indicate the best-performing method within each modality group.

169 Finally, we solve the following optimization problem for each timestep of the robot skill segment:

$$q_t^* = \arg \min_q \underbrace{\sum_{i=1}^N \| f_{\text{FK}}(q, {}^A k'_i) - ({}^W T_O(t) \cdot {}^O k'_i(t)) \|_2^2}_{\text{match keypoints}} + \underbrace{\lambda \| q - q_{t-1}^* \|_2^2}_{\text{temporal smoothness penalty}} \quad (2)$$

170 where $\lambda \geq 0$ trades off keypoint matching and temporal smoothness in joint space. This optimization
171 encourages a robot trajectory where the transformed actor keypoints ${}^A k'_i$ track the updated keypoint-
172 trajectory ${}^O k'_i(t)$ in the world frame under new object geometries and object poses. We solve the
173 optimization problem using the L-BFGS-B [41] gradient-based optimizer from SciPy [42].

174 4.3 Motion Planning

175 After obtaining the optimized robot configurations $Q^* = [q_1^*, \dots, q_H^*]$ from the skill segment adapta-
176 tion described above, we plan a collision-free trajectory from the robot’s current joint configura-
177 tion to the first joint configuration of the upcoming skill segment q_1^* . We query a collision-free motion
178 planner [43, 44] with the robot’s start configuration, the first joint configuration of the upcoming
179 skill, and provide collision geometries extracted directly from our simulation environment. As our
180 action space uses end-effector poses, we use an inverse kinematics controller to track the generated
181 trajectory.

182 5 Experiments

183 Through our experiments, we wish to answer the following questions:

184 **Q1:** Can a policy trained on data generated from a *single demonstration* using CP-Gen generalize
185 to *unseen object geometries and poses*?

186 **Q2:** How does CP-Gen’s *keypoint-trajectory formulation* affect *data generation* success rates on
187 tasks with *novel object geometries*?

188 **Q3:** How does training policies on *diverse object geometries* affect performance in various settings?

189 **Q4:** Can *pixel-based policies* trained with CP-Gen in simulation *transfer zero-shot* to the real world?

190 5.1 Simulation Experiments

191 **Tasks.** We evaluate all methods on single-arm Franka Panda tasks from the MimicGen [6] bench-
192 mark. We evaluate under two environment reset distributions. The default, fixed geometry distri-
193 bution varies only object poses, and corresponds to distribution “D1” in the MimicGen benchmark.
194 The custom *Geometry Generalization* distribution introduces shape variations via non-uniform scal-
195 ing and retains the same pose reset range as the default distribution (Figure 3).

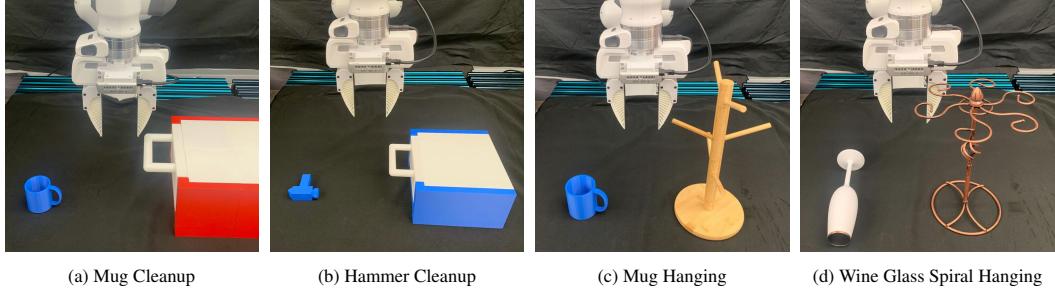


Figure 4: **Real World Tasks.** We evaluate on four challenging real-world tasks and show that policies trained on CP-Gen generated simulation datasets can transfer zero-shot to the real world.

Evaluation Metrics. We report the average success rate achieved by a particular policy network on the given task, along with the standard error of this metric. For computing standard errors, we use the best checkpoint rolled out for three seeds of 50 trials each.

Baselines. We compare CP-Gen to MimicGen [6], a data generation method that produces demonstrations in which object poses differ from those seen in source demonstrations. MimicGen first parses each source demonstration into object-centric subtask segments. To generate a new demonstration, it selects a reference segment, transforms it to match the new scene’s object pose, runs linear interpolation to reach the start of the transformed reference segment and executes the resulting end-effector trajectory using a controller. Finally, MimicGen filters out failed runs using a custom success checker, in the same way as CP-Gen. The authors additionally add Gaussian action noise to increase state coverage at the cost of the number of successfully generated demonstrations. In our case, we use 0.01 noise by default, but set noise to 0 if the number of successfully generated demonstrations drops below 10% for a given task. We use MimicGen to generate 1000 successful demonstrations per task from 10 human demonstrations.

Policy Training. For each method and task, we train a Diffusion Policy [45], specifically the DDPM variant, on the corresponding generated dataset that contains 1000 successful demonstrations. We use the same success classifier for both CP-Gen and MimicGen to filter for successful demonstrations. For policy inputs, we use two types of modalities from the third-person view and hand camera: (1) RGB images or (2) depth map and segmentation masks. We also provide proprioception information (end-effector pose and gripper width).

Results. From Table 1, we see that policies trained with data generated by CP-Gen (which takes a single source demonstration) outperform those trained using data generated using MimicGen across both observation input modalities. On the default task variants, we see that CP-Gen outperforms policies trained using MimicGen data (85% vs. 63%). We hypothesize that the policy performance drop is due to dataset quality, and hypothesize two reasons. First, CP-Gen uses a collision-free motion planner, and thus contains collision-avoidance behaviors, while MimicGen data does not have explicit collision-avoidance behaviors. Second, CP-Gen uses *actor keypoints* (keypoints on the robot or grasped object) during data generation which enables the robot to adapt actions based on how an object is grasped, while MimicGen only preserves the relative SE(3) transformation between the robot and a reference object. For example, on the Threading task, regardless of how the robot grasps the needle, MimicGen would generate the same end-effector actions for the robot in the threading phase based on the pose of the tripod object. In contrast, CP-Gen would adjust the robot’s actions to align the actor keypoints on the needle tip with the tripod. On the *Geometry Generalization* task variants, CP-Gen’s advantage is even more pronounced (70% vs. 37%). This finding is unsurprising, since CP-Gen generated demonstrations are adapted to the novel object geometries, while MimicGen does not consider changes to object geometry when generating demonstrations.

Method	Mug Cleanup	Hammer Cleanup	Mug Hanging	Wine Glass Spiral Hanging	Average
CP-Gen	0.80 ± 0.13	0.80 ± 0.13	1.00 ± 0.00	0.70 ± 0.15	0.83
MimicGen	0.20 ± 0.13	0.00 ± 0.00	0.80 ± 0.13	0.60 ± 0.16	0.40

Table 2: **CP-Gen policies successfully transfer zero-shot sim2real.** Policies take in depth and segmentation masks, and tasks feature both object-geometry and object pose variations.

(a) Data Generation			(b) Policy Evaluation		
Method	Default	Geometry	Method	Default	Geometry
MimicGen	0.58	0.23	CP-Gen (Default)	0.88	0.45
CP-Gen	0.89	0.62	CP-Gen (Geometry)	0.72	0.73

Table 3: **Ablations.** (a) CP-Gen outperforms MimicGen in terms of data generation success rate under both fixed (*Default*) and varied object geometry (*Geometry Generalization*) settings. (b) Using CP-Gen to generate diverse object geometries, denoted by *CP-Gen (Geometry)*, substantially improves policy generalization to novel geometries (73% vs. 45%). Furthermore, given a fixed training dataset budget of 1000 demonstrations, training and evaluating on the same distribution yields the best policy evaluation performance.

232 5.2 Real World Experiments

233 **Tasks.** We evaluate CP-Gen and baselines on four real-world manipulation tasks (Figure 4): (1)
234 *Mug Cleanup*: Open a drawer, pick and place mug into a drawer, and close the drawer. (2) *Hammer*
235 *Cleanup*: Same as *Mug Cleanup*, except with a small hammer, introducing grasping challenges.
236 Cleanup tasks test multi-stage, non-prehensile manipulation involving articulated objects. (3) *Mug*
237 *Hanging*: Pick and insert a mug onto a branch-like hook. This task tests constrained insertion
238 through a small opening. (4) *Wine Glass Spiral Hanging*: Pick, reorient and hang a wine glass by
239 threading the stem through a spiral-shaped rack. This task tests tight tolerance path following and
240 orientation control. For each task, we manually construct a digital twin in simulation, generate 1000
241 successful demonstrations using a single source demonstration, and train a Diffusion Policy on the
242 dataset. Policies are evaluated on two novel object geometries for each task.

243 **Policy Inputs.** To reduce the sim-to-real gap [46], we use depth maps and segmentation masks as
244 input observations. Specifically, we use depth and segmentation masks for the third-person-view
245 camera, and segmentation masks for the hand camera.

246 **Results.** From Table 2, we see that CP-Gen achieves strong zero-shot sim-to-real transfer on all four
247 manipulation tasks. In contrast, MimicGen struggles with novel geometries. In the *Cleanup* tasks,
248 failures arise from misalignment during drawer opening. In the *Hanging* tasks, failures are caused
249 by suboptimal grasp positioning, which prevents proper insertion.

250 5.3 Ablations

251 We conduct a series of ablations to study two aspects of CP-Gen: the use of the keypoint-trajectory
252 constraint formulation, and the effect of training on diverse object geometries. First, to assess the
253 impact of the keypoint-trajectory constraint formulation, we compare data generation success rates
254 between CP-Gen and MimicGen (which uses a relative pose data generation formulation) across
255 both default and *Geometry Generalization* task variants on eight simulation tasks. For each task, we
256 run 50 data generation trials. As shown in Table 3(a), CP-Gen achieves consistently higher success
257 rates (89% vs. 58%, 62% vs. 23%), especially on the *Geometry Generalization* reset distribution,
258 highlighting the value of the geometry-aware data generation that our keypoint-trajectory constraint
259 formulation enables. Second, we evaluate how training on diverse object geometries affects policy
260 generalization. In Table 3(b), we report policy success rates when training with and without geom-
261 etry diversity, and evaluated in the *Default* and *Geometry Generalization* environment reset settings.
262 As expected, training with diverse object geometries substantially improves policy generalization
263 to novel geometries (73% vs. 45%). We also observe that training and evaluating on the same
264 distribution yields the best evaluation performance (73% vs. 45% and 88% vs. 72%).

265 6 Conclusion

266 In this work, we propose CP-Gen, a data generation framework for generalizable visual imitation
267 learning given just one expert demonstration. CP-Gen generates geometry-aware demonstrations
268 using the insight that robot skills can be formulated as keypoint-trajectory constraints: keypoints on
269 the robot or grasped object must track a reference trajectory defined relative to a task-relevant object.
270 We demonstrate that CP-Gen’s approach to data generation from just one source demonstration
271 achieves state-of-the-art success rates on the MimicGen simulation benchmark, and outperforms
272 baselines on a custom benchmark featuring diverse object geometries. Finally, we show successful
273 zero-shot sim-to-real transfer of policies trained with CP-Gen data to challenging real-world tasks.

274 **Limitations.** CP-Gen requires manual annotation of robot skill segments and keypoints. The sim-
 275 ulation environment and reward functions are also manually defined; future work could incorporate
 276 foundation model based success classifiers. The method assumes a fixed skill sequence, limiting task
 277 level generalization. Incorporating task level planning via foundation models [47–50] or task and
 278 motion planner methods [51–53] can further improve data generation success. CP-Gen is less ef-
 279 fective for tasks with complex constraint sequences, such as dexterous in-hand manipulation, where
 280 methods such as reinforcement learning or model predictive control may be more suitable. Current
 281 results are limited to a single-arm Panda robot, though the method can be extended to multi-arm,
 282 mobile base, and legged systems [8]. CP-Gen does not handle the situation where one may want
 283 demonstrations generated for a specific object instance within the same category for which a source
 284 demonstration is available. For example, we may have a source demonstration for a mug with a square
 285 handle, but wish to generate demonstrations for a specific mug that might have a round han-
 286 dle. CP-Gen would need require the geometric transformation from the square handle mug to the
 287 round handle mug, which may be non-trivial to obtain. Future work may explore the use of neural
 288 descriptor fields [25] to assist in this setting.

289 References

- 290 [1] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn, “Learning Fine-Grained Bimanual Manipu-
 291 lation with Low-Cost Hardware,” in *Proceedings of Robotics: Science and Systems*, Daegu,
 292 Republic of Korea, Jul. 2023. DOI: [10.15607/RSS.2023.XIX.016](https://doi.org/10.15607/RSS.2023.XIX.016).
- 293 [2] Z. Fu, T. Z. Zhao, and C. Finn, “Mobile aloha: Learning bimanual mobile manipulation with
 294 low-cost whole-body teleoperation,” in *Conference on Robot Learning (CoRL)*, 2024.
- 295 [3] T. Z. Zhao *et al.*, “Aloha unleashed: A simple recipe for robot dexterity,” *arXiv preprint arXiv:2410.13126*, 2024.
- 297 [4] J. Aldaco *et al.*, “Aloha 2: An enhanced low-cost hardware for bimanual teleoperation,” *arXiv preprint arXiv:2405.02292*, 2024.
- 299 [5] A. Khazatsky *et al.*, “Droid: A large-scale in-the-wild robot manipulation dataset,” 2024.
- 300 [6] A. Mandlekar *et al.*, “Mimicgen: A data generation system for scalable robot learning using
 301 human demonstrations,” in *7th Annual Conference on Robot Learning*, 2023.
- 302 [7] C. R. Garrett, A. Mandlekar, B. Wen, and D. Fox, “Skillgen: Automated demonstration gener-
 303 ation for efficient skill learning and deployment,” in *8th Annual Conference on Robot Learn-
 304 ing*, 2024. [Online]. Available: <https://openreview.net/forum?id=YOFrRTDC6d>.
- 305 [8] Z. Jiang *et al.*, “Dexmimicgen: Automated data generation for bimanual dexterous manipu-
 306 lation via imitation learning,” in *2025 IEEE International Conference on Robotics and Au-
 307 tomation (ICRA)*, 2025.
- 308 [9] E. Ameperosa, J. A. Collins, M. Jain, and A. Garg, “Rocoda: Counterfactual data augmenta-
 309 tion for data-efficient robot learning from demonstrations,” *arXiv preprint arXiv:2411.16959*,
 310 2024.
- 311 [10] Z. Xue, S. Deng, Z. Chen, Y. Wang, Z. Yuan, and H. Xu, “Demogen: Synthetic demonstra-
 312 tion generation for data-efficient visuomotor policy learning,” *arXiv preprint arXiv:2502.16932*,
 313 2025.
- 314 [11] D. Wang, R. Walters, and R. Platt, “SO(2)-Equivariant Reinforcement Learning,” in *Inter-
 315 national Conference on Learning Representations*, 2022.
- 316 [12] C. Kohler, A. S. Srikanth, E. Arora, and R. Platt, “Symmetric models for visual force policy
 317 learning,” *arXiv preprint arXiv:2308.14670*, 2023.
- 318 [13] H. H. Nguyen, A. Baisero, D. Klee, D. Wang, R. Platt, and C. Amato, “Equivariant reinforce-
 319 ment learning under partial observability,” in *Conference on Robot Learning*, PMLR, 2023,
 320 pp. 3309–3320.
- 321 [14] H. Nguyen, T. Kozuno, C. C. Beltran-Hernandez, and M. Hamaya, “Symmetry-aware rein-
 322 forcement learning for robotic assembly under partial observability with a soft wrist,” *arXiv preprint arXiv:2402.18002*, 2024.

- 324 [15] B. Eisner, Y. Yang, T. Davchev, M. Vecerik, J. Scholz, and D. Held, “Deep SE(3)-equivariant
 325 geometric reasoning for precise placement tasks,” in *The Twelfth International Conference*
 326 *on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=2inBuwTyL2>.
- 328 [16] D. Wang, R. Walters, X. Zhu, and R. Platt, “Equivariant Q Learning in Spatial Action Spaces,”
 329 in *5th Annual Conference on Robot Learning*, 2021.
- 330 [17] H. Huang, D. Wang, A. Tangri, R. Walters, and R. Platt, “Leveraging Symmetries in Pick and
 331 Place,” *The International Journal of Robotics Research*, 2023.
- 332 [18] A. Simeonov *et al.*, “SE(3)-Equivariant Relational Rearrangement with Neural Descriptor
 333 Fields,” in *Conference on Robot Learning*, PMLR, 2023, pp. 835–846.
- 334 [19] C. Pan, B. Okorn, H. Zhang, B. Eisner, and D. Held, “TAX-Pose: Task-Specific Cross-
 335 Pose Estimation for Robot Manipulation,” in *Conference on Robot Learning*, PMLR, 2023,
 336 pp. 1783–1792.
- 337 [20] H. Huang, D. Wang, X. Zhu, R. Walters, and R. Platt, “Edge Grasp Network: A Graph-Based
 338 SE(3)-invariant Approach to Grasp Detection,” in *International Conference on Robotics and
 339 Automation (ICRA)*, 2023.
- 340 [21] S. Liu *et al.*, “Continual Vision-based Reinforcement Learning with Group Symmetries,” in
 341 *Conference on Robot Learning*, PMLR, 2023, pp. 222–240.
- 342 [22] M. Jia *et al.*, “Seil: Simulation-augmented equivariant imitation learning,” in *2023 IEEE In-
 343 ternational Conference on Robotics and Automation (ICRA)*, IEEE, 2023, pp. 1845–1851.
- 344 [23] S. Kim, B. Lim, Y. Lee, and F. C. Park, “Se (2)-equivariant pushing dynamics models for
 345 tabletop object manipulations,” in *Conference on Robot Learning*, PMLR, 2023, pp. 427–
 346 436.
- 347 [24] H. Huang, D. Wang, R. Walters, and R. Platt, “Equivariant Transporter Network,” in *Robotics:
 348 Science and Systems*, 2022.
- 349 [25] A. Simeonov *et al.*, “Neural descriptor fields: Se (3)-equivariant object representations for
 350 manipulation,” in *ICRA*, 2022.
- 351 [26] H. Ryu, H.-i. Lee, J.-H. Lee, and J. Choi, “Equivariant Descriptor Fields: SE(3)-Equivariant
 352 Energy-Based Models for End-to-End Visual Robotic Manipulation Learning,” in *The
 353 Eleventh International Conference on Learning Representations*, 2023.
- 354 [27] H. Huang, O. Howell, X. Zhu, D. Wang, R. Walters, and R. Platt, “Fourier Transporter: Bi-
 355 Equivariant Robotic Manipulation in 3D,” *arXiv preprint arXiv:2401.12046*, 2024.
- 356 [28] X. Zhu, D. Wang, O. Biza, G. Su, R. Walters, and R. Platt, “Sample Efficient Grasp Learning
 357 Using Equivariant Models,” in *Robotics: Science and Systems*, 2022.
- 358 [29] X. Zhu, D. Wang, G. Su, O. Biza, R. Walters, and R. Platt, “On robot grasp learning using
 359 equivariant models,” *Autonomous Robots*, vol. 47, no. 8, pp. 1175–1193, 2023.
- 360 [30] D. Wang *et al.*, “Equivariant diffusion policy,” in *8th Annual Conference on Robot Learning*,
 361 2024. [Online]. Available: <https://openreview.net/forum?id=WD2kUVLT1g>.
- 362 [31] J. Yang, Z.-a. Cao, C. Deng, R. Antonova, S. Song, and J. Bohg, “Equibot: Sim(3)-equivariant
 363 diffusion policy for generalizable and data efficient learning,” in *8th Annual Conference on
 364 Robot Learning*, 2024.
- 365 [32] J. Yang, C. Deng, J. Wu, R. Antonova, L. Guibas, and J. Bohg, “Equivact: Sim(3)-equivariant
 366 visuomotor policies beyond rigid object manipulation,” in *2024 IEEE International Confer-
 367 ence on Robotics and Automation (ICRA)*, IEEE, 2024, pp. 9249–9255.
- 368 [33] R. Hoque, A. Mandlekar, C. Garrett, K. Goldberg, and D. Fox, *Intervengen: Interventional
 369 data generation for robust and data-efficient robot imitation learning*, 2024. arXiv: 2405.
 370 01472 [cs.RO].
- 371 [34] W. Huang, C. Wang, Y. Li, R. Zhang, and L. Fei-Fei, “Rekep: Spatio-temporal reasoning of
 372 relational keypoint constraints for robotic manipulation,” in *8th Annual Conference on Robot
 373 Learning*, 2024. [Online]. Available: <https://openreview.net/forum?id=9iG3SEbMnL>.
- 374 [35] M. Oquab *et al.*, “Dinov2: Learning robust visual features without supervision,” *arXiv
 375 preprint arXiv:2304.07193*, 2023.

- 376 [36] OpenAI, *Gpt-4 technical report*, 2023. arXiv: 2303.08774 [cs.CL].
- 377 [37] L. Manuelli, W. Gao, P. Florence, and R. Tedrake, “Kpam: Keypoint affordances for category-
378 level robotic manipulation,” in *ISRR*, 2019.
- 379 [38] Z. Qin, K. Fang, Y. Zhu, F.-F. Li, and S. Savarese, “Keto: Learning keypoint representations
380 for tool manipulation,” in *Proceedings of the IEEE International Conference on Robotics and
381 Automation (ICRA)*, IEEE, 2020, pp. 7278–7285.
- 382 [39] P. Sundaresan *et al.*, “Learning rope manipulation policies using dense object descriptors
383 trained on synthetic depth data,” in *ICRA*, 2020.
- 384 [40] L. Manuelli, Y. Li, P. Florence, and R. Tedrake, “Keypoints into the future: Self-supervised
385 correspondence in model-based reinforcement learning,” *arXiv preprint arXiv:2009.05085*,
386 2020.
- 387 [41] R. H. Byrd, P. Lu, J. Nocedal, and C. Zhu, “A limited memory algorithm for bound con-
388 strained optimization,” *SIAM Journal on Scientific Computing*, vol. 16, no. 5, pp. 1190–1208,
389 1995. DOI: 10.1137/0916069. eprint: <https://doi.org/10.1137/0916069>. [Online]. Available:
390 <https://doi.org/10.1137/0916069>.
- 391 [42] P. Virtanen *et al.*, “SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python,”
392 *Nature Methods*, vol. 17, pp. 261–272, 2020. DOI: 10.1038/s41592-019-0686-2.
- 393 [43] B. Sundaralingam *et al.*, *Curobo: Parallelized collision-free minimum-jerk robot motion gen-
394 eration*, 2023. arXiv: 2310.17274 [cs.RO].
- 395 [44] W. Thomason, Z. Kingston, and L. E. Kavraki, “Motions in microseconds via vectorized
396 sampling-based planning,” in *IEEE International Conference on Robotics and Automation*,
397 2024. [Online]. Available: <http://arxiv.org/abs/2309.14545>.
- 398 [45] C. Chi *et al.*, “Diffusion policy: Visuomotor policy learning via action diffusion,” in *Proceed-
399 ings of Robotics: Science and Systems (RSS)*, 2023.
- 400 [46] M. Dalal *et al.*, “Local policies enable zero-shot long-horizon manipulation,” *International
401 Conference of Robotics and Automation*, 2025.
- 402 [47] M. Ahn *et al.*, “Do as i can, not as i say: Grounding language in robotic affordances,” *arXiv
403 preprint arXiv:2204.01691*, 2022.
- 404 [48] W. Huang *et al.*, “Inner monologue: Embodied reasoning through planning with language
405 models,” *arXiv preprint arXiv:2207.05608*, 2022.
- 406 [49] K. Lin, C. Agia, T. Migimatsu, M. Pavone, and J. Bohg, “Text2motion: From natural language
407 instructions to feasible plans,” *Autonomous Robots*, vol. 47, no. 8, pp. 1345–1365, 2023.
- 408 [50] B. Liu *et al.*, “Llm+ p: Empowering large language models with optimal planning profi-
409 ciency,” *arXiv preprint arXiv:2304.11477*, 2023.
- 410 [51] L. P. Kaelbling and T. Lozano-Pérez, “Hierarchical task and motion planning in the now,”
411 in *2011 IEEE International Conference on Robotics and Automation*, 2011, pp. 1470–1477.
412 DOI: 10.1109/ICRA.2011.5980391.
- 413 [52] C. R. Garrett, T. Lozano-Pérez, and L. P. Kaelbling, “Pddlstream: Integrating symbolic plan-
414 ners and blackbox samplers via optimistic adaptive planning,” in *Proceedings of the Interna-
415 tional Conference on Automated Planning and Scheduling*, vol. 30, 2020, pp. 440–448.
- 416 [53] C. R. Garrett *et al.*, “Integrated task and motion planning,” *Annual review of control, robotics,
417 and autonomous systems*, vol. 4, pp. 265–293, 2021.