

AXIS: A Growable Community-Driven Data Engine for Scalable Robot Manipulation

Anonymous Author(s)
Affiliation
Address
email

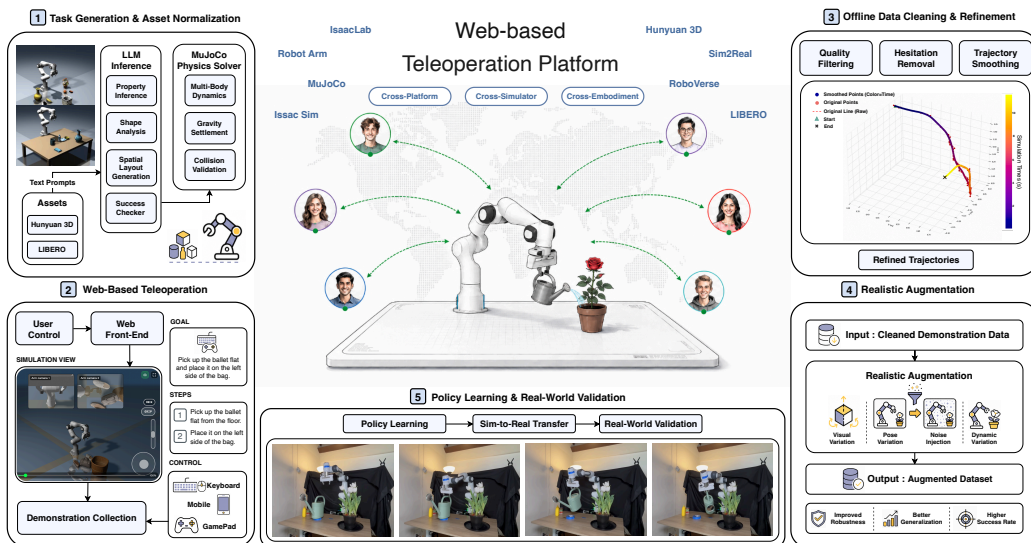


Figure 1: AXIS: A unified infrastructure bridging web-based teleoperation, GPU-accelerated realistic augmentation, and sim-to-real deployment.

1 **Abstract:** Robot manipulation policies require diverse, high-quality demonstra-
 2 tions, but existing data pipelines are often difficult to scale because they rely on
 3 specialized hardware, centralized operators, or fixed task suites. We present AXIS,
 4 a growable community-driven data engine and benchmark for scalable manipula-
 5 tion learning. AXIS collects demonstrations through browser-based MuJoCo-
 6 WASM teleoperation, generates and validates new manipulation tasks automati-
 7 cally, and converts community demonstrations into training-ready data through
 8 success checking, filtering, smoothing, resampling, and IsaacSim-based visual
 9 and physics augmentation. To make growth measurable, AXIS organizes data into
 10 task snapshots and evaluates policies with a fixed held-out protocol. We compare
 11 representative vision-language-action baselines under the same AXIS evaluation
 12 suite, and study scaling behavior across AXIS-100%, AXIS-50%, and AXIS-25%.
 13 These experiments test whether accessible community-collected data can support
 14 downstream policy learning, and whether increasing task coverage improves held-
 15 out manipulation generalization.

16 **Keywords:** Robot Learning, Manipulation Dataset, Data Collection

17 1 Introduction

18 General-purpose robot manipulation requires policies that can act across diverse objects, scenes, and
 19 task goals. Recent progress in imitation learning and vision-language-action modeling has shown
 20 that policy performance depends strongly on the scale, diversity, and quality of robot demonstra-
 21 tion data[1, 2]. However, collecting such data remains difficult. Physical robot demonstrations are

22 expensive, teleoperation systems often require specialized hardware or local simulator installations,
23 and many released manipulation datasets remain fixed after a one-time collection effort[3, 4]. As a
24 result, robot learning datasets often grow much more slowly than the models they are expected to
25 train.

26 A central limitation of current manipulation data pipelines is that data collection is often treated as
27 a closed and centralized process. Expert operators or small data teams collect demonstrations on
28 local hardware, process them offline, and release the resulting dataset as a static benchmark[4, 5].
29 This paradigm provides control over data quality, but it does not naturally support continuous task
30 expansion, broad community participation, or rapid iteration as models improve. Large-scale robot
31 learning instead requires data infrastructure that can keep growing along multiple axes: more tasks,
32 more objects, more scenes, more contributors, more visual conditions, and more physical variations.

33 We argue that the next generation of robot manipulation datasets should be *growable*: not only fixed
34 collections of trajectories, but mechanisms for continuously producing, validating, processing, and
35 evaluating new data. Growth should also be governed, so that new tasks and demonstrations are or-
36 ganized into reproducible snapshots and measured under a fixed evaluation protocol. To this end, we
37 introduce AXIS, a community-driven data engine and benchmark for scalable manipulation learning.
38 AXIS is designed around a simple principle: make data collection broadly accessible while scaling
39 high-quality data processing on demand. The system focuses on a Franka Research 3 arm with a
40 parallel-jaw gripper and supports tabletop manipulation tasks including pick-and-place, stacking,
41 pushing, pouring, articulated-object interaction, and tool-use-style manipulation. Each task includes
42 a language instruction, parameterized scene configuration, and structured success checker.

43 AXIS consists of three connected layers. The infrastructure layer combines automated task gener-
44 ation with browser-based MuJoCo-WASM teleoperation: new tasks are generated from high-level
45 instructions, object assets, scene layouts, and task-specific success conditions, then deployed to a
46 shared web interface where users collect demonstrations with commodity input devices. The dataset
47 layer converts raw community demonstrations into training-ready robot data by storing trajectories
48 in a unified format and applying success validation, filtering, static-segment removal, smoothing, re-
49 sampling, and IsaacSim-based visual and physics augmentation. The model layer uses the resulting
50 dataset to train and evaluate conventional visuomotor imitation learning and vision-language-action
51 policies under shared task definitions and success checkers. A key goal of AXIS is to make dataset
52 growth measurable. AXIS organizes training data into progressively larger task snapshots (AXIS-
53 25%, AXIS-50%, and AXIS-100%) and evaluates policies with a fixed held-out protocol, enabling
54 controlled scaling studies where the policy architecture, training budget, rollout protocol, evaluation
55 tasks, and success checkers remain fixed while only the training snapshot changes. Thus, we can
56 study whether increasing task coverage improves held-out manipulation generalization.

57 We evaluate AXIS through two primary experimental protocols. First, we compare representative
58 imitation learning and vision-language-action baselines under the same AXIS setup, testing whether
59 the dataset provides a consistent benchmark for model comparison. Second, we study dataset scaling
60 by training a fixed policy recipe across snapshots while keeping the evaluation suite and optimiza-
61 tion budget fixed. Together, these experiments examine whether community-collected, processed,
62 and augmented demonstrations can support downstream policy learning, and whether the growable
63 design of AXIS provides a practical path toward scalable robot manipulation.

64 **Contributions.** Our contributions are threefold. First, we present AXIS, a web-based and
65 community-driven infrastructure for scalable Franka manipulation data collection through MuJoCo-
66 WASM teleoperation and automated task generation. Second, we construct a growable manipulation
67 dataset with unified trajectory formatting, success validation, filtering, temporal smoothing, resam-
68 pling, IsaacSim-based visual and physics augmentation, and versioned task snapshots. Third, we in-
69 troduce an evaluation protocol for studying policy learning and within-embodiment dataset scaling,
70 including model baseline comparison and controlled AXIS-25%/50%/100% snapshot experiments
71 under a fixed held-out evaluation suite.

72 2 Related Work

73 **Large-scale and crowdsourced robot manipulation data.** Large-scale robot learning increasingly
74 relies on broad manipulation corpora spanning multi-robot transfer, cross-environment trajectories,
75 in-the-wild demonstrations, multi-embodiment data, and institutional aggregation [6, 7, 8, 9, 10,
76 11]. These datasets have supported scalable robot policies [1, 12], while recent benchmarks extend
77 coverage to lifelong learning, household tasks, compositional reasoning, memory, grasping, and
78 part-level supervision [13, 14, 15, 16, 17, 18, 19, 20]. Crowdsourced systems broaden collection [21,
79 22], but many resources remain fixed releases or depend on physical robots and specific setups.
80 AXIS is complementary: it uses browser-based simulated teleoperation for community collection
81 and organizes data into versioned snapshots, making dataset growth measurable rather than one-
82 time release.

83 **Teleoperation interfaces for robot data collection.** Robot demonstrations are commonly collected
84 by teleoperation, ranging from kinesthetic teaching [23] and leader–follower rigs [24, 25, 26] to
85 commodity devices [14, 27, 28, 29, 30, 31] and wearable capture [32, 33]. These interfaces trade
86 off precision, cost, embodiment fidelity, calibration burden, and accessibility. Web systems reduce
87 onboarding, but prior work often targets remote physical robots [22, 34], leaving collection con-
88 strained by robot availability, maintenance, and safety; high-fidelity simulators meanwhile require
89 heavier installation and GPU resources [35, 36]. AXIS separates collection from expensive sim-
90 ulation by using MuJoCo-WASM for browser teleoperation and reserving high-fidelity replay and
91 augmentation for backend processing.

92 **Simulation, task generation, and synthetic augmentation.** Simulation supports scalable task
93 variation, replay, and validation, with lightweight engines for interaction and prototyping [37,
94 38, 39, 40, 41] and GPU-oriented platforms for parallelism, rendering, and scene complex-
95 ity [42, 43, 44, 45, 35, 36]. Standardized task suites define common observations and success crite-
96 ria [46, 47, 48, 15, 2]; data-generation methods expand demonstrations by retargeting or synthesizing
97 trajectories [49, 50]; and LLM/VLM-driven systems scale task or asset diversity [51, 52, 53, 54].
98 AXIS connects task generation, browser teleoperation, and IsaacSim augmentation within a single
99 growable data pipeline.

100 **Vision-language-action models and manipulation evaluation.** Scaling data is closely tied to
101 training and evaluating general manipulation policies. Diffusion-style policies and action-chunking
102 transformers remain strong baselines, with ACT commonly used for chunked action prediction [25].
103 VLA models fine-tune large vision-language backbones on robot trajectories and have advanced
104 instruction-following manipulation [55, 56, 57, 58, 59, 60, 61, 62, 63, 64]. Evaluation has likewise
105 moved toward shared interfaces and robustness probes, including lifelong tasks, perturbation-aware
106 extensions, sim-to-real-correlated evaluation, and combinatorial scene shifts [13, 65, 66, 67]. AXIS
107 is complementary: rather than freezing a single suite, it releases versioned training snapshots with
108 a fixed held-out protocol so policies can be compared under shared data, observation, and success-
109 checker conventions as data grows.

110 3 Scalable Robot Data Collection

111 The dataset layer converts raw community demonstrations into training-ready robot trajectories
112 through validation, filtering, refinement, and augmentation. Since demonstrations are collected from
113 human operators rather than scripted controllers, the resulting data contains behavioral diversity in
114 approach direction, grasp timing, correction behavior, and execution style.

115 Completed demonstrations are serialized into a unified trajectory format. Each trajectory contains
116 task metadata, robot embodiment, environment identifier, simulator version, timestamps, observa-
117 tions, robot states, actions, and success information. Optional fields include object states, camera
118 streams, segmentation markers, and failure labels when available. This format makes demonstra-
119 tions replayable, filterable, and compatible with downstream policy learning.

120 Raw community demonstrations may contain jitter, stalls, corrupted frames, discontinuities, or
121 incorrect success labels. The dataset is therefore filtered using consistency checks for missing

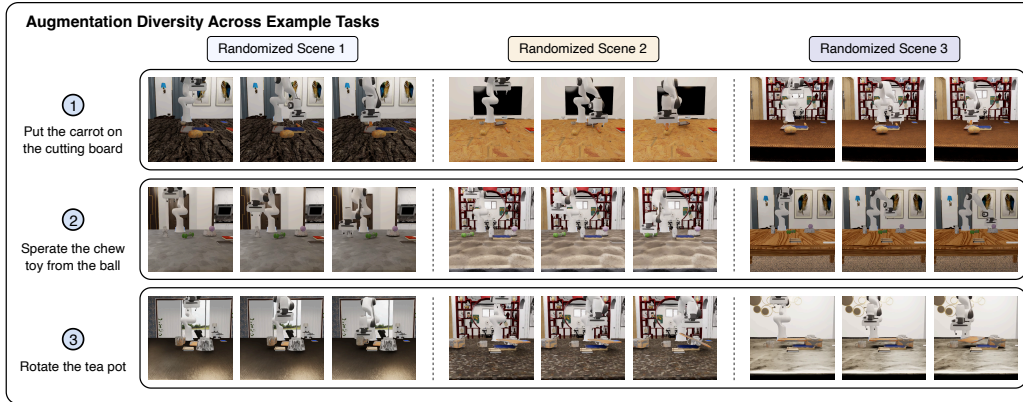


Figure 2: Visualization of augmentation diversity for three example tasks. For each task (row), we show three augmented variations (columns) obtained by randomizing appearance, lighting, environment, and initial object positions. Across augmentations, we vary materials and textures (e.g., table, floor, objects), lighting conditions (e.g., brightness and shadows), background and scene layout (e.g., walls, furniture, decorations), and initial states (e.g., object positions and orientations). This process substantially increases data diversity and improves policy robustness to environmental variations.

122 keys, anomalous numerical values, invalid state transitions, physically implausible frame-to-frame
 123 changes, and failed task completion. Although frontend dumps include success metadata, AXIS
 124 reuses task-specific success checkers during backend validation rather than relying solely on fron-
 125 tend success flags. Retained trajectories are then refined by removing static or corrupted segments,
 126 smoothing high-frequency teleoperation artifacts, and resampling to a fixed control frequency. This
 127 produces more stable action sequences while preserving the geometric structure of the demonstrated
 128 manipulation behavior.

129 Cleaned demonstrations are replayed in an IsaacSim backend for visual and physics augmenta-
 130 tion. Physical randomization varies object poses, clutter configurations, mass, friction, and related
 131 dynamics parameters. Rendering randomization varies camera viewpoints, lighting, textures, and
 132 visual appearance. These perturbations are applied while preserving task semantics, expanding the
 133 training distribution without requiring additional human demonstrations for every variation. The re-
 134 sult is a scalable pipeline that converts accessible web-collected demonstrations into higher-quality
 135 data for VLA and robot policy learning.

136 4 The AXIS Franka Dataset

137 In this section, we present the AXIS Franka Dataset and its associated validation, augmentation, and
 138 benchmarking framework.

139 **Dataset overview.** AXIS is a large-scale tabletop manipulation dataset built around a Franka Re-
 140 search 3 robot equipped with a parallel-jaw gripper. As summarized in Fig. 3, the current dataset
 141 snapshot contains 207 manipulation tasks, over 50k human demonstration trajectories, and more
 142 than 60k task or scene variants. Each task is specified by a language instruction, a parameterized
 143 simulation scene, task assets, and a structured success checker. Each demonstration trajectory con-
 144 tains synchronized multi-modal information, including RGB-D renderings from two camera views,
 145 robot states, object states, robot actions, task metadata, and success labels. This unified data schema
 146 enables the same trajectories to be used for imitation learning, replay-based verification, visual aug-
 147 mentation, and benchmark evaluation. The task suite covers a broad range of common tabletop
 148 manipulation behaviors, including pick-and-place, relocation, stacking, sorting, pushing, pouring,
 149 insertion, alignment, and interaction with articulated or container-like objects, and the current task
 150 distribution spans seven scene categories.

151 **Validation and refinement.** To ensure data quality, all demonstrations undergo a standardized val-
 152 idation pipeline that verifies trajectory completeness, task success, and physical consistency. Valid
 153 trajectories are subsequently smoothed and resampled to a unified control frequency, improving

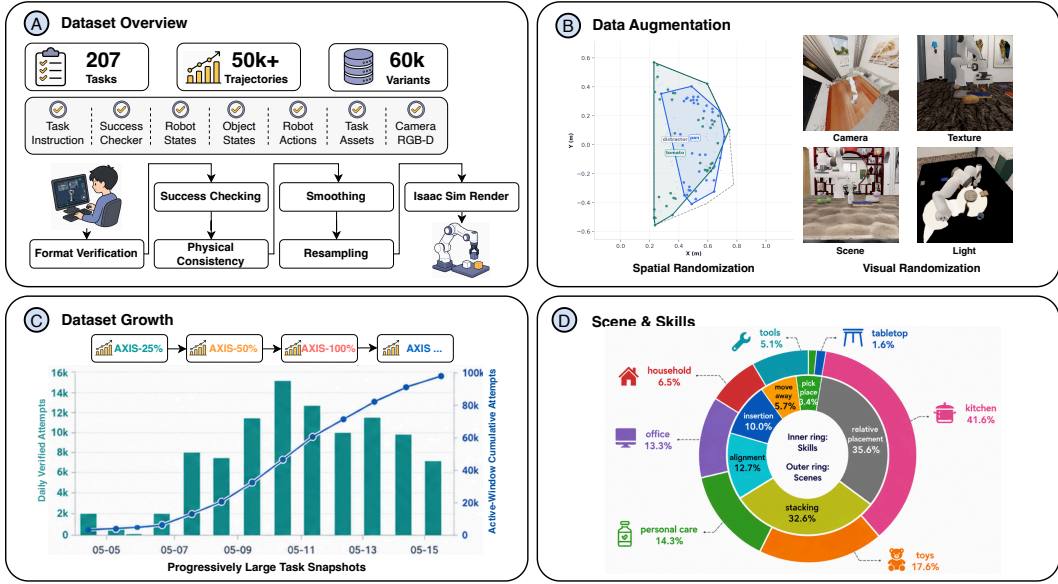


Figure 3: Overview of AXIS datasets.

Table 1: Effect of data refinement on teleoperation trajectory quality. Relative to raw teleoperation, the final smoothed-and-resampled trajectories reduce mean acceleration and mean jerk by 63.9% and 80.8%, respectively.

Data Version	Sampling Rate	Mean Acceleration	Mean Jerk	Replay Success
Raw Teleoperation	5.0 Hz	1.3539	11.5899	100.0%
Smoothed	5.0 Hz	0.6382	2.9160	91.4%
Smoothed + Resampled	20 Hz	0.4885	2.2243	86.2%

154 temporal consistency and reducing teleoperation noise. As shown in Table 1, the refinement process
 155 substantially reduces acceleration and jerk while maintaining a high replay success rate. Additional
 156 implementation details are provided in the Appendix ??.

157 **Data augmentation.** To improve diversity and robustness, AXIS incorporates augmentation through
 158 spatial and visual randomization. These augmentations generate additional valid scene configura-
 159 tions while preserving task semantics and success conditions. As illustrated in Fig. 2, the resulting
 160 variations cover changes in object layouts, viewpoints, textures, lighting, and backgrounds, expand-
 161 ing the dataset to more than 60k scene variants. This diversity makes AXIS particularly useful for
 162 studying robustness and generalization under environmental variations.

163 **Benchmarking platform.** AXIS is designed as a growable benchmark for studying policy learn-
 164 ing as data grows. The training data is organized into progressively larger task snapshots, including
 165 AXIS-25%, AXIS-50%, and AXIS-100%. Each snapshot preserves the same data format, validation
 166 pipeline, task definitions, rollout budget, and success checkers, enabling controlled scaling studies
 167 in which only the amount and diversity of training data are varied. The dataset is also continuously
 168 expandable through community-driven task and trajectory collection. In the active collection win-
 169 dows shown in Fig. 3, the number of daily verified attempts reaches approximately 15k at its peak,
 170 and the active-window cumulative number of attempts approaches 100k. These statistics indicate
 171 that AXIS is intended to support incremental dataset growth rather than a one-time static release.

172 The benchmark supports two complementary evaluation modes. First, AXIS enables direct compar-
 173 ison across representative policy families, including conventional visuomotor imitation learning
 174 methods and vision-language-action policies. Second, AXIS supports within-embodiment scaling
 175 studies, where policies are trained on AXIS-25%, AXIS-50%, AXIS-100%, and future larger snap-
 176 shots to measure whether additional task coverage improves held-out manipulation generalization.

177 Together, these properties position AXIS not only as a large-scale manipulation dataset, but also as
178 a long-term benchmark for studying scalable robot learning and data-driven generalization.

179 5 Experiments

180 To evaluate AXIS as a growable data engine, we organize our experiments around three ques-
181 tions. (i) Does pretraining $\pi_{0.5}$ [62] on AXIS improve downstream LIBERO-Plus [65] performance,
182 and is the improvement specific to AXIS versus a matched-volume baseline drawn from Robo-
183 Casa365 [15]? (ii) Does the improvement scale with AXIS volume, and is it largest when the down-
184 stream LIBERO fine-tuning data is scarce? (iii) On which LIBERO-Plus perturbation axes does
185 AXIS pretraining help most, and does the pattern match the perturbations explicitly randomized by
186 the AXIS IsaacSim augmentation pipeline?

187 5.1 Experimental Setup

188 **AXIS dataset.** The AXIS training corpus is collected through the browser-based MuJoCo-WASM
189 teleoperation frontend and stored in the unified AXIS trajectory format. The robot embodiment is
190 a Franka Research 3 arm with a parallel-jaw gripper; tasks cover tabletop manipulation behaviors
191 including pick-and-place, stacking and sorting, articulated object interaction, pushing, pouring, and
192 tool-use-style manipulation. Raw demonstrations are passed through the AXIS processing pipeline
193 (success validation, static-segment removal, Savitzky–Golay smoothing, fixed-rate resampling) and
194 then replayed in IsaacSim under randomized visual and physical conditions. We use the AXIS-100%
195 snapshot as the largest pretraining pool, and sample uniformly without replacement at the task level
196 to construct the 25% and 50% subsets used in scaling experiments.

197 **Training pipeline.** All conditions share the same three-stage training pipeline: start from the re-
198 leased $\pi_{0.5}$ checkpoint; optionally pretrain on a sim corpus (AXIS subset or RoboCasa365 subset);
199 fine-tune on the LIBERO trajectory set provided with LIBERO-Plus. The pretraining stage uses
200 identical optimizer settings, learning rate, batch size, and gradient-step budget across all non-vanilla
201 conditions; the fine-tuning stage uses the $\pi_{0.5}$ default LIBERO recipe across all conditions. The
202 RoboCasa-matched control is constructed by sampling RoboCasa365 trajectories without replace-
203 ment until the trajectory count equals AXIS-100%, controlling for raw pretraining volume.

204 **Evaluation protocol.** All policies are evaluated on the LIBERO-Plus robustness suite, which ex-
205 tends LIBERO with systematic perturbations along seven axes (Camera, Light, Sensor Noise, Back-
206 ground, Layout, Language, Robot). For each (task, perturbation axis) pair we fix the rollout budget
207 to a constant K across conditions, and report mean success rate over the rollouts. As a within-
208 distribution sanity check, we additionally report success rate on the held-out task split of AXIS-
209 100% itself; this confirms that AXIS pretraining does not destructively shift the policy away from
210 the AXIS task distribution.

211 5.2 Main Result: AXIS Pretraining vs. Matched Sim Baseline

212 Our first experiment tests whether AXIS data is useful as pretraining for a downstream manipulation
213 benchmark, and whether the benefit is specific to AXIS or generic to any matched sim corpus. We
214 treat AXIS as *pretraining* data on top of $\pi_{0.5}$ [62], fine-tune the resulting model on the standard
215 LIBERO trajectories, and evaluate on the LIBERO-Plus robustness benchmark [65]. This design
216 holds the downstream task family, the fine-tuning data, and the evaluation distribution fixed for
217 every method, isolating the effect of the pretraining corpus.

218 We compare four pretraining conditions. (1) $\pi_{0.5}$ **vanilla** performs the LIBERO fine-tune directly
219 from the released $\pi_{0.5}$ checkpoint and establishes the floor. (2) $\pi_{0.5}$ + **AXIS-25%** pretrains on a 25%
220 random subset of AXIS-100% before LIBERO fine-tuning, providing a low-volume scaling point.
221 (3) $\pi_{0.5}$ + **AXIS-100%** pretrains on the full AXIS-100% snapshot and is our headline condition.
222 (4) $\pi_{0.5}$ + **RoboCasa-matched** pretrains on a RoboCasa365 [15] subset matched in trajectory count
223 to AXIS-100%, controlling for the possibility that any comparably-sized Franka-sim pretraining
224 corpus would yield similar gains. Together, conditions (1)–(3) measure whether AXIS helps and

Table 2: Main result. All conditions start from the released $\pi_{0.5}$ checkpoint, optionally pretrain on the listed corpus, then fine-tune on LIBERO trajectories using identical hyperparameters, and evaluate on LIBERO-Plus. The RoboCasa-matched control uses the same trajectory count as AXIS-100%. Per-perturbation columns report success rate on the corresponding LIBERO-Plus axis.

Pretraining	# pretrain demos	Overall \uparrow	Cam.	Light	Noise	B.G.	Layout	Lang.	Robot
$\pi_{0.5}$ vanilla	0	66.5	50.3	90.6	61.8	87.7	61.7	71.9	54.8
$\pi_{0.5}$ + AXIS-25%	0.25 N_{AXIS}	72.0	56.1	91.8	68.1	89.1	72.1	76.0	61.3
$\pi_{0.5}$ + AXIS-50%	0.50 N_{AXIS}	75.1	59.7	92.5	72.5	89.6	77.9	78.3	65.1
$\pi_{0.5}$ + AXIS-100% (ours)	N_{AXIS}	79.4	65.5	93.1	78.5	90.4	84.9	81.1	70.2
$\pi_{0.5}$ + RoboCasa-matched (ctrl.)	N_{AXIS}	57.5	35.2	79.5	63.2	81.7	68.0	49.2	39.4

225 whether more AXIS helps more; condition (4) measures whether the gain is AXIS-specific or generic
 226 to sim pretraining at this volume.

227 All conditions share the same downstream LIBERO fine-tuning recipe ($\pi_{0.5}$ default optimizer, batch
 228 size, schedule, and step budget), the same LIBERO trajectory set, the same LIBERO-Plus evaluation
 229 protocol, and the same number of rollouts per task and perturbation condition. Pretraining hyperpa-
 230 rameters are also matched across conditions (2)–(4): identical learning rate, batch size, gradient-step
 231 budget, and optimizer state initialization.

232 The key question addressed by Table 2 is whether AXIS pretraining provides benefits beyond the
 233 existing pretraining of $\pi_{0.5}$, and whether those benefits are specific to AXIS data. AXIS pretraining
 234 consistently improves downstream performance over the vanilla $\pi_{0.5}$ baseline, with gains increasing
 235 as more AXIS data is added. The performance gap between AXIS-25% and AXIS-100% further
 236 suggests that these benefits have not yet saturated, indicating that additional AXIS data continues to
 237 improve robustness on downstream manipulation tasks. The matched-volume RoboCasa365 control
 238 tests whether these gains arise from AXIS itself or simply from additional simulation data. Despite
 239 using the same pretraining volume, RoboCasa-matched underperforms both AXIS variants across
 240 nearly all perturbation categories and on the overall benchmark. This suggests that the gains are not
 241 explained by dataset size alone, but by properties unique to the AXIS pipeline, including crowd-
 242 sourced task creation, diverse demonstrations, and large-scale environment randomization.

243 Overall, the results support AXIS as a scalable data engine: increasing AXIS data yields consis-
 244 tent downstream improvements, and those gains stem from the diversity and coverage of the AXIS
 245 pipeline rather than pretraining volume alone.

246 5.3 Per-Perturbation Robustness Analysis

247 To understand where AXIS pretraining helps and where it does not, we report the per-perturbation
 248 breakdown of the AXIS-100% condition against vanilla $\pi_{0.5}$ and the RoboCasa-matched control at
 249 the full LIBERO fine-tune budget. The LIBERO-Plus perturbation axes probe different distribution
 250 shifts: *Camera* varies viewpoint, *Light* varies illumination, *Sensor Noise* adds photometric pertur-
 251 bations (Gaussian noise and contrast jitter) to the camera observations, *Background* varies scene
 252 appearance and table/wall textures, *Layout* varies object initial positions and orientations, *Language*
 253 varies task wording, and *Robot* varies the robot initial pose. Different pretraining corpora are ex-
 254 pected to help unequally across these axes depending on what their training distribution covers;
 255 AXIS’s IsaacSim augmentation specifically targets camera, light, texture, and layout perturbations,
 256 so we expect the strongest relative gains on those axes.

257 We use Table 3 to understand where AXIS pretraining helps and where it does not. As expected from
 258 the AXIS augmentation design, the largest gains occur on visual and scene-related perturbations, in-
 259 cluding camera viewpoint, sensor noise, background appearance, and object layout, showing that
 260 IsaacSim randomization transfers effectively to corresponding robustness challenges in LIBERO-
 261 Plus. However, AXIS also improves performance on language and robot-pose variation, despite
 262 not explicitly augmenting these factors. This suggests that AXIS provides more than perturbation-
 263 specific robustness: exposure to diverse tasks, behaviors, and interaction contexts improves the
 264 underlying manipulation representations, leading to broader gains under distribution shift. Compar-
 265 isons with the RoboCasa365 baseline of the same volume reinforce this conclusion. AXIS outper-

Table 3: Per-perturbation LIBERO-Plus breakdown at the full LIBERO fine-tune budget. Each row reports success rate on the corresponding perturbation axis. Δ_{van} is the absolute improvement of AXIS-100% over $\pi_{0.5}$ vanilla. AXIS-RC is the difference between AXIS-100% and the matched-volume RoboCasa365 control, measuring the gain over a same-size sim corpus rather than per-trajectory superiority.

Perturbation axis	$\pi_{0.5}$ vanilla	+ AXIS-25%	+ AXIS-100%	+ RoboCasa-m.	Δ_{van}	AXIS-RC
Camera	50.3	56.1	65.5	35.2	+15.2	+30.3
Light	90.6	91.8	93.1	79.5	+2.5	+13.6
Sensor Noise	61.8	68.1	78.5	63.2	+16.7	+15.3
Background	87.7	89.1	90.4	81.7	+2.7	+8.7
Layout	61.7	72.1	84.9	68.0	+23.2	+16.9
Language	71.9	76.0	81.1	49.2	+9.2	+31.9
Robot	54.8	61.3	70.2	39.4	+15.4	+30.8
Overall	66.5	72.0	79.4	57.5	+12.9	+21.9

forms the control across nearly all perturbation axes, indicating that its benefits stem not only from additional simulation data but also from the diversity of environments, tasks, and demonstrations it provides. Overall, AXIS improves robustness broadly, with particularly strong gains on the visual and geometric shifts it was designed to address.

6 Discussion

Why growable datasets matter. Static datasets are valuable for benchmarking, but robot manipulation learning also benefits from data sources that can grow with model capability and observed failures. AXIS turns data collection into an iterative feedback loop: TaskGen expands semantic diversity through new manipulation problems, community teleoperation adds behavioral diversity through varied human strategies and corrections, and visual and physics augmentation expands environmental diversity through changes in appearance, scene configuration, and dynamics. In this sense, dataset growth becomes a mechanism for improving policy robustness rather than a one-time release process.

Community collection and data quality. Community-driven collection improves scale and accessibility, but introduces variation in operator skill, motion smoothness, strategy, and task completion quality. Useful variation exposes policies to diverse grasps, recoveries, and correction patterns, while harmful variation appears as corrupted states, idle segments, failed rollouts, discontinuities, or unstable actions. AXIS addresses this tension by making validation, cleaning, smoothing, replay, and augmentation core stages of dataset construction, combining broad participation with the quality control needed for downstream policy learning.

7 Conclusion

We introduced AXIS, a growable community-driven data engine and benchmark for scalable Franka manipulation learning. AXIS combines browser-based MuJoCo-WASM teleoperation, automated task generation, distributed demonstration collection, standardized trajectory processing, IsaacSim-based augmentation, and fixed-protocol evaluation. It constructs a training-ready dataset that can expand through new tasks, demonstrations, and augmented conditions, while enabling controlled comparisons across vision-language-action policies. More broadly, AXIS points toward continuously extensible robot data pipelines that connect task generation, community collection, policy training, and failure-driven improvement.

Limitations and Future Work AXIS currently focuses on simulated Franka tabletop manipulation, and sim-to-real transfer remains an open challenge. Future work will expand AXIS to more robot embodiments, richer sensing modalities, longer-horizon tasks, finer-grained annotations, and active failure-driven data collection to support more general and scalable robot learning.

299 **References**

- 300 [1] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan,
301 K. Hausman, et al. Rt-1: Robotics transformer for real-world control at scale, 2023. URL
302 <https://arxiv.org/abs/2212.06817>.
- 303 [2] H. Geng, F. Wang, S. Wei, Y. Li, B. Wang, B. An, C. T. Cheng, H. Lou, P. Li, Y.-J. Wang,
304 Y. Liang, D. Goetting, C. Xu, H. Chen, Y. Qian, Y. Geng, J. Mao, W. Wan, M. Zhang, J. Lyu,
305 S. Zhao, J. Zhang, J. Zhang, C. Zhao, H. Lu, Y. Ding, R. Gong, Y. Wang, Y. Kuang, R. Wu,
306 B. Jia, C. Sferrazza, H. Dong, S. Huang, Y. Wang, J. Malik, and P. Abbeel. Roboverse: Towards
307 a unified platform, dataset and benchmark for scalable and generalizable robot learning, 2025.
308 URL <https://arxiv.org/abs/2504.18904>.
- 309 [3] J. Duan, S. Yu, H. L. Tan, H. Zhu, and C. Tan. A survey of embodied ai: From simulators to
310 research tasks. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(2):
311 230–244, 2022.
- 312 [4] H. Choi, C. Crump, C. Duriez, A. Elmquist, G. Hager, D. Han, F. Hearl, J. Hodgins, A. Jain,
313 F. Leve, et al. On the use of simulation in robotics: Opportunities, challenges, and suggestions
314 for moving forward. *Proceedings of the National Academy of Sciences*, 118(1):e1907856118,
315 2021.
- 316 [5] J. Eßer, N. Bach, C. Jestel, O. Urbann, and S. Kerner. Guided reinforcement learning: A
317 review and evaluation for efficient and effective real-world robotics [survey]. *IEEE Robotics
318 & Automation Magazine*, 30(2):67–85, 2022.
- 319 [6] S. Dasari, F. Ebert, S. Tian, S. Nair, B. Bucher, K. Schmeckpeper, S. Singh, S. Levine, and
320 C. Finn. Robonet: Large-scale multi-robot learning, 2020. URL [https://arxiv.org/abs/
321 1910.11215](https://arxiv.org/abs/1910.11215).
- 322 [7] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and
323 S. Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets,
324 2021. URL <https://arxiv.org/abs/2109.13396>.
- 325 [8] H. R. Walke, K. Black, T. Z. Zhao, Q. Vuong, C. Zheng, P. Hansen-Estruch, et al. Bridge-
326 data v2: A dataset for robot learning at scale. In J. Tan, M. Toussaint, and K. Darvish,
327 editors, *Proceedings of The 7th Conference on Robot Learning*, volume 229 of *Proceed-
328 ings of Machine Learning Research*, pages 1723–1736. PMLR, 06–09 Nov 2023. URL
329 <https://proceedings.mlr.press/v229/walke23a.html>.
- 330 [9] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, et al. Droid: A large-scale in-the-wild robot
331 manipulation dataset, 2025. URL <https://arxiv.org/abs/2403.12945>.
- 332 [10] K. Wu, C. Hou, J. Liu, Z. Che, X. Ju, Z. Yang, M. Li, Y. Zhao, Z. Xu, G. Yang, et al. Robo-
333 mind: Benchmark on multi-embodiment intelligence normative data for robot manipulation.
334 In *Robotics: Science and Systems (RSS) 2025*. Robotics: Science and Systems Foundation,
335 2025. URL <https://www.roboticsproceedings.org/rss21/p152.pdf>.
- 336 [11] A. O’Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, et al. Open X-Embodiment:
337 Robotic learning datasets and RT-X models. In *2024 IEEE International Conference on
338 Robotics and Automation (ICRA)*, pages 6892–6903, 2024. doi:10.1109/ICRA57147.2024.
339 10611477.
- 340 [12] B. Zitkovich, T. Yu, S. Xu, P. Xu, T. Xiao, F. Xia, J. Wu, et al. Rt-2: Vision-language-action
341 models transfer web knowledge to robotic control. In J. Tan, M. Toussaint, and K. Darvish,
342 editors, *Proceedings of The 7th Conference on Robot Learning*, volume 229 of *Proceedings
343 of Machine Learning Research*, pages 2165–2183. PMLR, 06–09 Nov 2023. URL <https://proceedings.mlr.press/v229/zitkovich23a.html>.

- 345 [13] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone. Libero: Bench-
346 marking knowledge transfer for lifelong robot learning. In A. Oh, T. Naumann,
347 A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural In-*
348 *formation Processing Systems*, volume 36, pages 44776–44791. Curran Associates, Inc.,
349 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/8c3c666820ea055a77726d66fc7d447f-Paper-Datasets_and_Benchmarks.pdf)
350 [8c3c666820ea055a77726d66fc7d447f-Paper-Datasets_and_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/8c3c666820ea055a77726d66fc7d447f-Paper-Datasets_and_Benchmarks.pdf).
- 351 [14] S. Nasiriany, A. Maddukuri, L. Zhang, A. Parikh, A. Lo, A. Joshi, A. Mandlekar, and Y. Zhu.
352 Robocasa: Large-scale simulation of everyday tasks for generalist robots. In *Robotics: Science*
353 *and Systems (RSS)*, 2024.
- 354 [15] S. Nasiriany, S. Nasiriany, A. Maddukuri, and Y. Zhu. Robocasa365: A large-scale simulation
355 framework for training and benchmarking generalist robots. In *International Conference on*
356 *Learning Representations (ICLR)*, 2026.
- 357 [16] Y. Chen, Z. Chen, N. T. Chan, J. Chen, J. Yin, J. Shi, Y. Gao, Y.-L. Li, and J. Huo. Robo-
358 himan: A hierarchical evaluation paradigm for compositional generalization in long-horizon
359 manipulation, 2025. URL <https://arxiv.org/abs/2510.13149>.
- 360 [17] Y. Dai, H. Fu, J. Lee, Y. Liu, H. Zhang, J. Yang, C. Finn, N. Fazeli, and J. Chai. Robomme:
361 Benchmarking and understanding memory for robotic generalist policies, 2026. URL <https://arxiv.org/abs/2603.04639>.
362
- 363 [18] S. K. Srinivas, Y. Shukla, A. Arnold, and S. Chitta. Graspfactory: A large object-centric
364 grasping dataset, 2025. URL <https://arxiv.org/abs/2509.20550>.
- 365 [19] P. Li, H. Geng, J. Crate, Y. Han, J. Zhang, F. Wang, C. T. Cheng, R. Dong, Y.-J. Wang, H. Lou,
366 et al. Rose: Reconstructing objects, scenes, and trajectories from casual videos for robotic
367 manipulation. In *NeurIPS 2025 Workshop on Bridging Language, Agent, and World Models*
368 *for Reasoning and Planning*.
- 369 [20] Y. Yin, Z. Han, S. Aarya, S. Xu, J. Wang, J. Peng, A. Wang, A. Yuille, and T. Shu. Partin-
370 struct: Part-level instruction following for fine-grained robot manipulation. In *Proceedings of*
371 *Robotics: Science and Systems (RSS)*, 2025.
- 372 [21] A. Mandlekar, Y. Zhu, A. Garg, J. Booher, M. Spero, A. Tung, J. Gao, J. Emmons, A. Gupta,
373 E. Orbay, S. Savarese, and L. Fei-Fei. RoboTurk: A crowdsourcing platform for robotic skill
374 learning through imitation. In *Conference on Robot Learning (CoRL)*, 2018.
- 375 [22] S. Mirchandani, M. Tang, J. Duan, J. I. Hamid, M. Cho, and D. Sadigh. Robocode: Gamify-
376 ing robot data collection. *arXiv preprint arXiv:2512.21235*, 2025. doi:10.48550/arXiv.2512.
377 21235.
- 378 [23] P. Kormushev, S. Calinon, and D. G. Caldwell. Imitation learning of positional and force skills
379 demonstrated via kinesthetic teaching and haptic input. *Advanced Robotics*, 25(5):581–603,
380 2011. doi:10.1163/016918611X558261.
- 381 [24] P. Wu, Y. Shentu, Z. Yi, X. Lin, and P. Abbeel. Gello: A general, low-cost, and intuitive
382 teleoperation framework for robot manipulators. *arXiv preprint arXiv:2309.13037*, 2023. doi:
383 10.48550/arXiv.2309.13037.
- 384 [25] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation
385 with low-cost hardware. In *Robotics: Science and Systems*, 2023. doi:10.15607/RSS.2023.
386 XIX.016.
- 387 [26] Y. Zou, Z. Zhou, C. Shi, Z. Ye, J. Huang, Y. Ding, and B. Zhao. U-arm: Ultra low-cost general
388 teleoperation interface for robot manipulation. *arXiv preprint arXiv:2509.02437*, 2025. doi:
389 10.48550/arXiv.2509.02437.

- 390 [27] D. Honerkamp, H. Mahesheka, J. O. von Hartz, T. Welschehold, and A. Valada. Zero-cost
391 whole-body teleoperation for mobile manipulation. *IEEE Robotics and Automation Letters*,
392 2025.
- 393 [28] C. Zhou, C. Peers, Y. Wan, R. Richardson, and D. Kanoulas. Teleman: Teleoperation for
394 legged robot loco-manipulation using wearable imu-based motion capture. *arXiv preprint*
395 *arXiv:2209.10314*, 2022. doi:10.48550/arXiv.2209.10314.
- 396 [29] A. George, A. Bartsch, and A. B. Farimani. Openvr: Teleoperation for manipulation. *Soft-*
397 *wareX*, 29:102054, 2025. doi:10.1016/j.softx.2025.102054.
- 398 [30] K. Lu, Y. He, C. Lu, and P. Li. I know kung fu: Synthetic dexterous hand demonstration
399 collection via vr teleoperation. In *NeurIPS 2025 Workshop on Space in Vision, Language, and*
400 *Embodied AI*.
- 401 [31] J. Wang, C.-C. Chang, J. Duan, D. Fox, and R. Krishna. Eve: Enabling anyone to train robot
402 using augmented reality. *arXiv preprint arXiv:2404.06089*, 2024. doi:10.48550/arXiv.2404.
403 06089.
- 404 [32] C. Wang, H. Shi, W. Wang, R. Zhang, L. Fei-Fei, and C. K. Liu. Dexcap: Scalable and portable
405 mocap data collection system for dexterous manipulation. *arXiv preprint arXiv:2403.07788*,
406 2024. doi:10.48550/arXiv.2403.07788.
- 407 [33] M. Xu, H. Zhang, Y. Hou, Z. Xu, L. Fan, M. Veloso, and S. Song. Dexumi: Using hu-
408 man hand as the universal manipulation interface for dexterous manipulation. *arXiv preprint*
409 *arXiv:2505.21864*, 2025. doi:10.48550/arXiv.2505.21864.
- 410 [34] C.-L. Fok, F. Sun, M. Mangum, A. K. Mok, B. He, and L. Sentis. Web-based teleoperation of
411 a humanoid robot. *arXiv preprint arXiv:1607.05402*, 2016. doi:10.48550/arXiv.1607.05402.
- 412 [35] NVIDIA. Isaac sim. <https://developer.nvidia.com/isaac/sim>, 2025. Robotics simu-
413 lator, accessed 2026-03-06.
- 414 [36] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, L. Yi,
415 A. X. Chang, L. J. Guibas, and H. Su. Sapien: A simulated part-based interactive environment.
416 *arXiv preprint arXiv:2003.08515*, 2020. doi:10.48550/arXiv.2003.08515.
- 417 [37] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In
418 *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–
419 5033. IEEE, 2012.
- 420 [38] N. Koenig and A. Howard. Design and use paradigms for gazebo, an open-source multi-robot
421 simulator. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages
422 2149–2154, Sendai, Japan, Sept. 2004. doi:10.1109/IROS.2004.1389727.
- 423 [39] E. Rohmer, S. P. N. Singh, and M. Freese. Coppeliassim (formerly v-rep): A versa-
424 tile and scalable robot simulation framework. In *2013 IEEE/RSJ International Confer-*
425 *ence on Intelligent Robots and Systems*, 2013. doi:10.1109/IROS.2013.6696520. Available:
426 www.coppeliarobotics.com.
- 427 [40] Robot Locomotion Group. Drake: Model-based design and verification for robotics. <https://drake.mit.edu/>, 2026. Accessed: 2026-03-06.
- 428 [41] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation for games, robotics
429 and machine learning. <http://pybullet.org>, 2016. Accessed: 2026-03-06.
- 430 [42] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin,
431 A. Allshire, A. Handa, and G. State. Isaac gym: High performance gpu-based physics simu-
432 lation for robot learning. In *Proceedings of the Neural Information Processing Systems Track*
433 *on Datasets and Benchmarks*, 2021. doi:10.48550/arXiv.2108.10470. NeurIPS 2021 Datasets
434 and Benchmarks Track.
435

- 436 [43] K. Zakka, B. Tabanpour, Q. Liao, M. Haiderbhai, S. Holt, J. Y. Luo, A. Allshire, E. Frey,
437 K. Sreenath, L. A. Kahrs, C. Sferrazza, Y. Tassa, and P. Abbeel. Mujoco playground: An open-
438 source framework for gpu-accelerated robot learning and sim-to-real transfer. *arXiv preprint*
439 *arXiv:2502.08844*, 2025. doi:10.48550/arXiv.2502.08844.
- 440 [44] C. D. Freeman, E. Frey, A. Raichuk, S. Girgin, I. Mordatch, and O. Bachem. Brax: A
441 differentiable physics engine for large scale rigid body simulation. In *Proceedings of the*
442 *Neural Information Processing Systems Track on Datasets and Benchmarks*, 2021. doi:
443 10.48550/arXiv.2106.13281. NeurIPS 2021 Datasets and Benchmarks Track.
- 444 [45] Genesis Authors. Genesis: A universal and generative physics engine for robotics and be-
445 yond. <https://github.com/Genesis-Embodied-AI/Genesis>, 2024. Software project,
446 accessed 2026-03-06.
- 447 [46] J. Gu, F. Xiang, X. Li, Z. Ling, X. Liu, T. Mu, Y. Tang, S. Tao, X. Wei, Y. Yao, X. Yuan,
448 P. Xie, Z. Huang, R. Chen, and H. Su. ManiSkill2: A unified benchmark for generalizable
449 manipulation skills. In *International Conference on Learning Representations (ICLR)*, 2023.
- 450 [47] S. Tao, F. Xiang, A. Shukla, Y. Qin, X. Hinrichsen, X. Yuan, C. Bao, X. Lin, Y. Liu, T. kai
451 Chan, Y. Gao, X. Li, T. Mu, N. Xiao, A. Gurha, Z. Huang, R. Calandra, R. Chen, S. Luo,
452 and H. Su. ManiSkill3: Gpu parallelized robotics simulation and rendering for generalizable
453 embodied AI. *arXiv preprint arXiv:2410.00425*, 2024.
- 454 [48] Y. Zhu, J. Wong, A. Mandlekar, R. Martín-Martín, A. Joshi, S. Nasiriany, and Y. Zhu. ro-
455 bosuite: A modular simulation framework and benchmark for robot learning. *arXiv preprint*
456 *arXiv:2009.12293*, 2020.
- 457 [49] A. Mandlekar, S. Nasiriany, B. Wen, I. Akinola, Y. Narang, L. Fan, Y. Zhu, and D. Fox.
458 MimicGen: A data generation system for scalable robot learning using human demonstrations.
459 In *Conference on Robot Learning (CoRL)*, 2023.
- 460 [50] Z. Jiang, Y. Xie, K. Lin, Z. Xu, W. Wan, A. Mandlekar, L. Fan, and Y. Zhu. DexMimicGen:
461 Automated data generation for bimanual dexterous manipulation via imitation learning. *arXiv*
462 *preprint arXiv:2410.24185*, 2024.
- 463 [51] L. Wang, Y. Ling, Z. Yuan, M. Shridhar, C. Bao, Y. Qin, B. Wang, H. Xu, and X. Wang.
464 GenSim: Generating robotic simulation tasks via large language models. In *International*
465 *Conference on Learning Representations (ICLR)*, 2024.
- 466 [52] P. Hua, M. Liu, A. Macaluso, Y. Lin, W. Zhang, H. Xu, and L. Wang. GenSim2: Scaling robot
467 data generation with multi-modal and reasoning LLMs. In *Conference on Robot Learning*
468 *(CoRL)*, 2024.
- 469 [53] Y. Wang, Z. Xian, F. Chen, T.-H. Wang, Y. Wang, K. Fragkiadaki, Z. Erickson, D. Held,
470 and C. Gan. RoboGen: Towards unleashing infinite data for automated robot learning via
471 generative simulation. In *International Conference on Machine Learning (ICML)*, 2024.
- 472 [54] Y. He, X. Wang, P. Li, Y. Huang, J. Masterjohn, J. Wu, L. Guibas, Y. Yang, Y. Jiang, and
473 C. Jiang. Fishbone: From one 3d asset to a million controllable edits, 2026. URL <https://arxiv.org/abs/2605.24805>.
474
- 475 [55] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna,
476 T. Kreiman, C. Xu, J. Luo, Y. L. Tan, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and
477 S. Levine. Octo: An open-source generalist robot policy. In *Robotics: Science and Systems*
478 *(RSS)*, 2024.
- 479 [56] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Fos-
480 ter, G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine,
481 P. Liang, and C. Finn. OpenVLA: An open-source vision-language-action model. In *Confer-*
482 *ence on Robot Learning (CoRL)*, 2024.

- 483 [57] M. J. Kim, C. Finn, and P. Liang. Fine-tuning vision-language-action models: Optimizing
484 speed and success. *arXiv preprint arXiv:2502.19645*, 2025.
- 485 [58] S. Liu, L. Wu, B. Li, H. Tan, H. Chen, Z. Wang, K. Xu, H. Su, and J. Zhu. RDT-1B: a diffusion
486 foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024.
- 487 [59] Q. Li, Y. Liang, Z. Wang, L. Luo, X. Chen, M. Liao, F. Wei, Y. Deng, S. Xu, Y. Zhang,
488 X. Wang, B. Liu, J. Fu, J. Bao, D. Chen, Y. Shi, J. Yang, and B. Guo. CogACT: A foundational
489 vision-language-action model for synergizing cognition and action in robotic manipulation.
490 *arXiv preprint arXiv:2411.19650*, 2024.
- 491 [60] H. Shi, B. Xie, Y. Liu, L. Sun, F. Liu, T. Wang, E. Zhou, H. Fan, X. Zhang, and G. Huang.
492 MemoryVLA: Perceptual-cognitive memory in vision-language-action models for robotic ma-
493 nipulation. *arXiv preprint*, 2025. arXiv.
- 494 [61] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Haus-
495 man, B. Ichter, S. Jakubczak, T. Jones, L. Ke, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair,
496 K. Pertsch, L. X. Shi, J. Tanner, Q. Vuong, A. Walling, H. Wang, and U. Zhilinsky. π_0 : A
497 vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*,
498 2024.
- 499 [62] Physical Intelligence, K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail,
500 M. Equi, C. Finn, N. Fusai, M. Y. Galliker, D. Ghosh, L. Groom, K. Hausman, B. Ichter,
501 S. Jakubczak, T. Jones, L. Ke, D. LeBlanc, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair,
502 K. Pertsch, A. Z. Ren, L. X. Shi, L. Smith, J. T. Springenberg, K. Stachowicz, J. Tanner,
503 Q. Vuong, H. Walke, A. Walling, H. Wang, L. Yu, and U. Zhilinsky. $\pi_{0.5}$: a vision-language-
504 action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.
- 505 [63] NVIDIA, J. Bjorck, F. C. neda, N. Cherniadev, X. Da, R. Ding, L. Fan, Y. Fang, D. Fox, F. Hu,
506 S. Huang, J. Jang, Z. Jiang, J. Kautz, K. Kundalia, L. Lao, Z. Li, Z. Lin, K. Lin, G. Liu,
507 E. Llontop, L. Magne, A. Mandlekar, A. Narayan, S. Nasiriany, S. Reed, Y. L. Tan, G. Wang,
508 Z. Wang, J. Wang, Q. Wang, J. Xiang, Y. Xie, Y. Xu, Z. Xu, S. Ye, Z. Yu, A. Zhang, H. Zhang,
509 Y. Zhang, Y. Zhang, L. Fan, and Y. Zhu. GR00T N1: An open foundation model for generalist
510 humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.
- 511 [64] B. Chen, T. Zhang, H. Geng, C. Zhang, P. Li, K. Song, W. T. Freeman, J. Malik, P. Abbeel,
512 R. Tedrake, et al. Large video planner enables generalizable robot control. *arXiv preprint*
513 *arXiv:2512.15840*, 2025.
- 514 [65] S. Fei, S. Wang, J. Shi, Z. Dai, J. Cai, P. Qian, L. Ji, X. He, S. Zhang, Z. Fei, J. Fu, J. Gong, and
515 X. Qiu. LIBERO-Plus: In-depth robustness analysis of vision-language-action models. *arXiv*
516 *preprint arXiv:2510.13626*, 2025.
- 517 [66] X. Li, K. Hsu, J. Gu, K. Pertsch, O. Mees, H. R. Walke, C. Fu, I. Lunawat, I. Sieh, S. Kir-
518 mani, S. Levine, J. Wu, C. Finn, H. Su, Q. Vuong, and T. Xiao. Evaluating real-world robot
519 manipulation policies in simulation. *arXiv preprint arXiv:2405.05941*, 2024.
- 520 [67] W. Pumacay, I. Singh, J. Duan, R. Krishna, J. Thomason, and D. Fox. THE COLOSSEUM:
521 A benchmark for evaluating generalization for robotic manipulation. In *Robotics: Science and*
522 *Systems (RSS)*, 2024.

523	Appendix Contents	
524	A Dataset Comparison	15
525	B Web-Infra: Browser-Based Infrastructure	16
526	C TaskGen: Task Generation and Deployment	17
527	D Data Cleaning and Trajectory Refinement	18
528	E AXIS Data Rendering	19
529	F $\pi_{0.5}$ Pretraining from Base	22
530	G LIBERO Post-training from AXIS-Pretrained Parameters	23
531	H Real-time Inference and Real-world Rollouts	24

532 **A Dataset Comparison**

Table 4: Comparison of representative robot manipulation datasets in terms of collection setting, scale, task coverage, and supervision.

Dataset	Year	Setting	Scale	Task Coverage	Data / Supervision
RoboNet [6]	2019	Real robot	162K trajectories	Diverse manipulation tasks	RGB videos, robot actions, gripper states, and action trajectories
LIBERO [13]	2023	Simulation	6.5K trajectories	130 tasks	RGB video, proprioception, language instructions, and expert demonstrations
BridgeData V2 [8]	2023	Real robot	60K trajectories; 24 environments	13 skills	RGB videos, proprioception, goal images, and language instructions
Open X-Embodiment [11]	2024	Real robot	1M+ episodes; 60 datasets	500+ skills	RGB videos, end-effector poses, language instructions, and action trajectories
DROID [9]	2024	Real robot	76K trajectories; 564 scenes	86 tasks	RGB-D videos, proprioception, teleoperated actions, and language instructions
RoboMIND [10]	2024	Real robot	107K trajectories; 96 objects	279 tasks	RGB-D videos, proprioception, language, teleoperation, and failure demonstrations
PartInstruct [20]	2025	Simulation	513 objects; 1,302 tasks; 10K+ demos	16 part-level task classes	Language, part-level 3D object data, skill sequences, and expert demonstrations
Ours	2026	Web-based simulation	50,129 episodes	207 tasks	Browser-based episodes with scalable user-contributed manipulation data

533 Table 4 summarizes representative robot manipulation datasets spanning real-world robotics, sim-
534 ulation environments, and large-scale embodied learning benchmarks. Existing datasets gener-
535 ally emphasize either large-scale real-world data collection (e.g., RoboNet, Open X-Embodiment,
536 DROID, and RoboMIND) or highly controlled simulation benchmarks with expert demonstrations
537 (e.g., LIBERO and PartInstruct). While these datasets have substantially advanced robot learning,
538 their collection pipelines often require specialized robotic hardware, expert operators, or centralized
539 infrastructure, limiting scalability and accessibility.

540 In contrast, our dataset adopts a web-based simulation paradigm that enables robot manipulation
541 data collection directly through standard web browsers. The framework emphasizes three key prop-
542 erties: scalability, continual growth, and community participation. The absence of physical robots,
543 specialized infrastructure, and expert operators significantly lowers participation barriers and en-
544 ables large-scale data contribution from a broad user community. New tasks, objects, and inter-
545 action scenarios can be incorporated continuously as the community expands. The current release
546 contains 50,129 episodes spanning 207 manipulation tasks, illustrating the potential of communi-
547 ty-driven data generation for robotic manipulation. We envision this framework as a complementary
548 direction to existing real-world and simulator-based datasets and as a sustainable pathway toward
549 ever-growing robot manipulation corpora.

550 **B Web-Infra: Browser-Based Infrastructure**

551 AXIS is a cross-simulator, cross-machine infrastructure for scalable robot demonstration collection,
 552 trajectory processing, and policy evaluation. The key system decision is to separate interactive data
 553 collection from compute-heavy backend processing. Contributors teleoperate robots in a browser-
 554 based MuJoCo WebAssembly frontend using commodity devices, while uploaded demonstrations
 555 are routed to backend machines for validation, cleaning, replay, rendering, augmentation, and ex-
 556 port to downstream policy-learning pipelines. This split lowers the entry cost for data contributors
 557 without forcing the full pipeline to inherit the limits of a lightweight browser simulator.

558 At the pipeline level, AXIS follows six stages: browser-based teleoperation and logging, trajectory
 559 upload in a unified format, backend replay and realistic augmentation, data cleaning and curation,
 560 model training with sim-to-real evaluation, and real-robot deployment. The pipeline is intentionally
 561 asymmetric: latency-sensitive interaction remains on commodity web clients, while throughput-
 562 sensitive rendering, replay, model training, and evaluation run on backend machines. Data render-
 563 ing is performed on 8×RTX 4090 GPUs, while all model training and evaluation are performed on
 564 8×A100 GPUs. This organization lets the same demonstration serve multiple purposes, including
 565 task verification, cleaned control supervision, realistic RGB rendering, domain-randomized aug-
 566 mentation, downstream policy learning, and evaluation.

567 The infrastructure is organized around shared task and trajectory abstractions. At the task level,
 568 environments expose common observation, action, reset, step, and success-checking interfaces, so a
 569 task collected in MuJoCo-WASM can later be verified, replayed, or rendered in a backend simulator
 570 without rewriting the task logic. At the trajectory level, each demonstration stores metadata, time-
 571 series signals, and episode-level annotations in a unified format. Metadata include the task name,
 572 environment identifier, robot embodiment, simulator version, contributor identifier when available,
 573 and timestamps. Time-series fields include actions, robot states, observations, and optional object
 574 states or visual observations. Episode-level fields include success labels, segmentation markers,
 575 failure modes when available, and links to the task-specific success checker.

Table 5: Functional decomposition of the Web-Infra pipeline.

Stage	Execution site	Role in the pipeline
Task deployment	Browser client	Loads a parameterized task, robot model, scene objects, and success checker into the MuJoCo-WASM frontend.
Teleoperation	Browser client	Maps keyboard, mouse, joystick, or gamepad commands to end-effector targets and joint-level robot actions.
Logging	Browser client	Buffers time-aligned state-action samples inside the physics loop and serializes successful attempts with task metadata.
Upload and indexing	Backend service	Stores trajectories with versioned metadata so they can be retrieved for validation, replay, rendering, and training.
Cleaning and replay	Backend machines	Filters invalid episodes, smooths and resamples motion, verifies success, and reconstructs trajectories for rendering.
Learning export	Backend machines	Packages cleaned and rendered data for imitation learning and VLA training.

576 The browser frontend is implemented as an end-to-end teleoperation interface that covers task se-
 577 lection, interactive 3D visualization, camera control, control-state display, and automatic trajectory
 578 packaging. The core dashboard acts as a heads-up display over the MuJoCo-WASM and WebGL
 579 viewport. A top-level status area reports real-time metrics such as elapsed time and rendering rate,
 580 while the control panel groups gripper motion, Cartesian movement, rotation, action buttons, and
 581 camera controls so new contributors can quickly understand the available inputs. The interface sup-
 582 ports keyboard, mouse, virtual joystick, and gamepad control, making data collection compatible
 583 with heterogeneous personal computers rather than specialized lab hardware.

584 To keep browser control responsive, AXIS uses an orchestrated runtime rather than a monolithic
 585 UI loop. MuJoCo physics stepping and Three.js rendering are separated from the main React UI

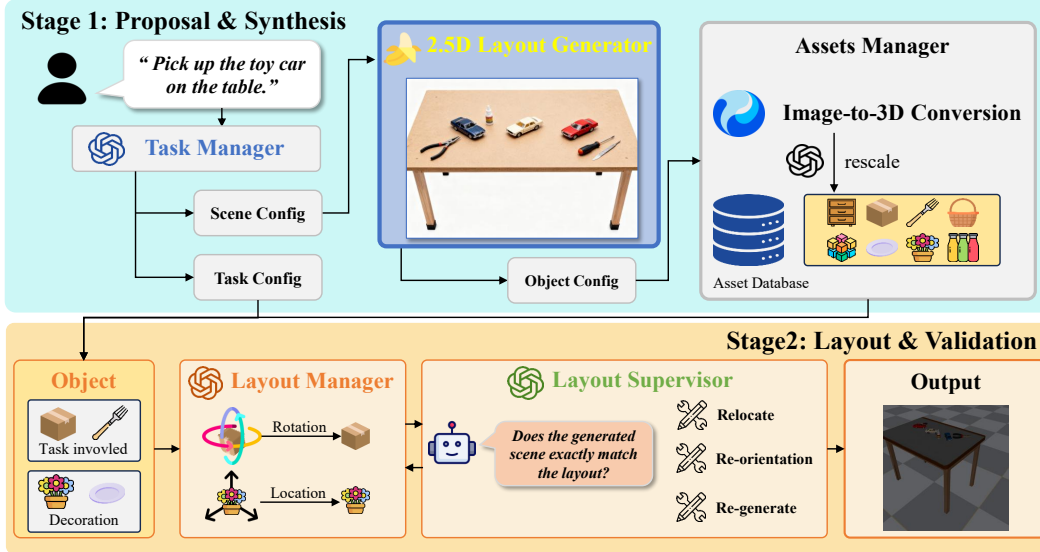


Figure 4: TaskGen pipeline overview. A language instruction is converted into task, scene, and object specifications; object models are retrieved or generated and normalized; a layout is proposed, instantiated, checked, and deployed.

586 thread. A lightweight broadcaster sends occasional status updates to the interface, while the tra-
 587 jectory manager buffers state-action samples directly inside the high-frequency control loop. This
 588 design avoids frequent UI re-rendering during teleoperation and keeps the logged samples aligned
 589 with the simulator state. Once the embedded checker verifies task completion, the buffered samples
 590 and task metadata are packaged into a standardized trajectory record for backend upload.

591 Operator commands are mapped to six-degree-of-freedom Cartesian end-effector targets and con-
 592 verted to joint-level robot commands through inverse kinematics. For more complex end-effectors,
 593 AXIS also supports interpolation-based grasp control. Let $\mathbf{q}_{\text{rest}} \in \mathbb{R}^n$ denote the fully open hand
 594 pose and $\mathbf{q}_{\text{tmpl}} \in \mathbb{R}^n$ denote a grasp template, such as a power, precision, or functional grasp. The
 595 commanded hand pose is

$$\mathbf{q}(t) = \mathbf{q}_{\text{rest}} + \alpha(t) (\mathbf{q}_{\text{tmpl}} - \mathbf{q}_{\text{rest}}), \quad (1)$$

596 where $\alpha(t) \in [0, 1]$ is controlled by the operator. This collapses a high-dimensional hand config-
 597 uration, often with more than 16 actuated joints, into a single scalar control axis. New grippers
 598 or dexterous hands can therefore be integrated by defining rest poses and task-specific template
 599 libraries, without changing the core teleoperation controller.

600 C TaskGen: Task Generation and Deployment

601 TaskGen converts a high-level manipulation instruction into a simulation-ready task instance. Given
 602 an instruction such as *pick up the toy car on the table*, a task manager decomposes the request into
 603 task, scene, and object configurations. The task configuration specifies the manipulation goal and
 604 required object interactions. The scene configuration defines workspace context such as supporting
 605 surfaces, clutter level, and spatial constraints. The object configuration determines which object
 606 models appear in the scene, their semantic roles, and whether they are required for task completion
 607 or included as decorative context. Each task is also assigned a difficulty level from 1 to 5, which
 608 controls the number of objects, the tightness of spatial constraints, and the overall scene complexity.

609 TaskGen uses an object-model manager to populate scenes with diverse objects. Object models can
 610 be retrieved from an existing database or generated through an image-to-3D conversion pipeline.
 611 Because generated meshes may have inconsistent physical scales, the object-model manager nor-
 612 malizes them to plausible dimensions before adding them to the task database. This normalization

613 step is important for downstream teleoperation and replay: object size affects reachability, grasp
614 feasibility, collision behavior, and the validity of geometric success conditions.

615 Conditioned on the task specification, selected object models, and difficulty level, a 2.5D layout gen-
616 erator proposes an initial scene arrangement. The generator describes which objects should appear
617 on the workspace and their approximate spatial relationships. A layout manager then instantiates the
618 proposal into a full 3D environment by assigning each object a precise location and orientation. A
619 layout supervisor verifies whether the instantiated scene satisfies the intended task constraints. When
620 the supervisor detects an inconsistency, such as an unreachable object, an invalid support relation,
621 or a mismatch between the language goal and object placement, the system iteratively corrects the
622 scene through relocation, re-orientation, or partial regeneration.

Table 6: TaskGen pipeline components and their outputs.

Component	Input	Output / validation role
Task manager	Language instruction, difficulty level	Structured task, scene, and object configurations with semantic object roles.
Object-model manager	Object requirements and model database	Retrieved or image-to-3D generated object models normalized to plausible physical scale.
2.5D layout generator	Task specification, object models, clutter settings	Initial object arrangement over the workspace, including task-involved and decorative objects.
Layout manager	2.5D proposal and object geometry	Full 3D scene with object poses, orientations, support relations, and workspace placement.
Layout supervisor	Instantiated scene and task constraints	Validated scene or corrective actions such as relocation, re-orientation, or partial regeneration.
Success checker	Task goal and object interactions	Automatic completion test reused during browser collection and backend verification.

623 After a valid layout is produced, TaskGen defines distributions over initial conditions to support sce-
624 nario sampling. These distributions vary object poses, clutter configurations, and camera viewpoints
625 while preserving the core task semantics. Each generated instance is paired with a task-specific
626 success checker that encodes completion using geometric constraints and object interactions. The
627 checker is used both during browser teleoperation, where it terminates and packages successful
628 demonstrations, and during backend validation, where it helps verify whether uploaded trajectories
629 remain task-consistent after transfer.

630 The output of TaskGen has three parts: a parameterized task definition describing scene configu-
631 ration and object layout, a unified observation and action interface for execution across simulators,
632 and a task-specific success checker for automatic evaluation. This decomposition makes task growth
633 continuous rather than one-off. New language instructions, objects, layouts, and difficulty settings
634 can be added to the task pool, deployed to the browser, collected by contributors, and then replayed
635 or rendered through the same backend pipeline.

636 **D Data Cleaning and Trajectory Refinement**

637 Raw web demonstrations require filtering and refinement before they can be used for policy learning.
638 Browser teleoperation can introduce jitter, stalled samples, unintended oscillations, low-frequency
639 sampling artifacts, corrupted records, and occasional task-completion mismatches. AXIS therefore
640 treats data cleaning as a production pipeline with two goals: filter unusable episodes and refine
641 borderline episodes into temporally smooth, replayable trajectories.

642 The validation stage removes trajectories that violate basic structural, numerical, physical, or task-
643 level checks. Structural checks reject demonstrations with missing required fields, incomplete
644 metadata, corrupted records, or length mismatches between state and action streams. Numerical
645 checks reject non-finite values and anomalous records. Physical checks compute successive
646 changes in robot states and reject trajectories whose deltas exceed predefined plausible thresholds,
647 which catches extreme discontinuities and invalid state transitions. Task-level checks remove failed

Algorithm 1 Trajectory cleaning and refinement.

Require: Raw episode S , static threshold ϵ , raw timestep Δt_{raw} , target timestep Δt , Savitzky-Golay parameters W, P

Ensure: Smoothed and resampled episode S_{final}

- 1: **Step 1: Validation and segmentation**
 - 2: Reject S if required fields, finite values, length consistency, physical continuity, or task success checks fail.
 - 3: $S_{\text{clean}} \leftarrow \text{REMOVESTATICSEGMENTS}(S, \epsilon)$
 - 4: $T_{\text{clean}} \leftarrow \text{REBUILDUNIFORMTIME}(S_{\text{clean}}, \Delta t_{\text{raw}})$
 - 5: $Q_{\text{clean}} \leftarrow \text{ROBOTJOINTSEQUENCE}(S_{\text{clean}})$
 - 6: **Step 2: Robot-motion smoothing**
 - 7: $Q_{\text{smooth}} \leftarrow \text{SAVITZKYGOLAY}(Q_{\text{clean}}, W, P)$
 ▷ Discrete gripper transitions can be excluded from continuous smoothing when needed.
 - 8: **Step 3: Fixed-frequency resampling**
 - 9: $T_{\text{new}} \leftarrow \text{ARRANGE}(T_{\text{clean}}[0], T_{\text{clean}}[-1], \Delta t)$
 - 10: $Q_{\text{final}} \leftarrow \text{CUBICSPLINE}(T_{\text{clean}}, Q_{\text{smooth}})(T_{\text{new}})$
 - 11: **Step 4: Extra-state alignment and export**
 - 12: $S_{\text{final}} \leftarrow \text{ALIGNNONROBOTSTATE}(T_{\text{new}}, Q_{\text{final}}, S_{\text{clean}})$ **return** S_{final}
-

648 episodes and can re-verify nominal success labels using the task-specific checker in the backend
649 simulation environment.

650 After validation, the pipeline removes uninformative static segments. Samples are classified as static
651 when the absolute variation across all robot joints is below a configured threshold; in the current
652 processing pipeline, this threshold is 5×10^{-3} . Removing these samples eliminates artificial pauses
653 and idle segments caused by human hesitation, while preserving the portions of the episode that
654 contain meaningful control changes. The remaining trajectory is then re-timed to form a consistent
655 temporal sequence for smoothing and resampling.

656 Continuous robot action or joint trajectories are smoothed with a Savitzky-Golay filter. The current
657 implementation uses a window size of 15 and a polynomial order of 3 for robot motion smoothing,
658 while discrete gripper transitions can be excluded from the smoothing pass so that open-close events
659 do not become physically ambiguous. This stage reduces high-frequency teleoperation artifacts
660 such as abrupt corrections and small oscillations without replacing the original demonstration with
661 an unrelated planned path.

662 Because browser demonstrations may be recorded at a low or non-uniform rate, cleaned trajectories
663 are resampled to a fixed control frequency with cubic spline interpolation. The web interface can
664 produce trajectories at only 6–8 Hz, while downstream control and policy training require a higher
665 frequency; the current target frequency is 20 Hz. Robot states and continuous action targets are
666 interpolated after smoothing. Non-robot states, such as object poses, bypass the action-smoothing
667 filter and are temporally aligned through interpolation only. The resulting trajectory remains faithful
668 to the demonstrated task execution while becoming smoother, denser, and better aligned with replay
669 and policy-learning requirements.

670 The cleaned output serves as the handoff point to later rendering and augmentation stages. A retained
671 episode contains validated metadata, temporally aligned robot and object states, smoothed robot
672 motion, success annotations, and episode boundaries. The same cleaned record can be replayed for
673 verification, rendered into RGB observation streams, perturbed through IsaacSim scene or camera
674 randomization, or exported for imitation-learning and VLA training. In this way, data cleaning is
675 not only a quality-control step but also the interface that makes web-collected demonstrations usable
676 by the larger AXIS dataset and model-training pipeline.

677 E AXIS Data Rendering

678 AXIS converts verified task trajectories into LIBERO-style visual imitation-learning data through
679 state replay. The rendering source is not an action log that must be re-executed through a simulator.
680 Instead, each eligible attempt is represented as packed simulator state, which provides the robot and

Table 7: Quantitative evaluation of trajectory optimization across 16 LIBERO tasks. Smoothness (Max Acceleration and Max Jerk) is compared before and after the proposed processing.

LIBERO Task	Smoothness				Pos. Dev. (m)	Removed Ratio
	Mean Acc.		Mean Jerk			
	Before	After	Before	After		
Task 1: Place Black Bowl on Top of cabinet	0.2458	0.0738	2.6425	1.3363	0.0235	6.46%
Task 2: Place Rear Butter in Cabinet Top Drawer and Close It	0.3067	0.0943	3.7395	1.6208	0.0654	5.73%
Task 3: Place the black bowl on the plate	0.3559	0.1011	5.1358	2.3634	0.0201	7.17%
Task 4: Place the black bowl on top of the cabinet	0.3225	0.0859	4.6724	2.0785	0.0142	3.77%
Task 5: Place the frying pan on the stove	0.1832	0.1055	2.4771	1.1374	0.0106	1.24%
Task 6: Place the moka pot on the stove	0.1819	0.1090	2.5239	1.6784	0.0064	1.94%
Task 7: Turn on the stove	0.2689	0.1295	3.8291	2.3017	0.0142	3.81%
Task 8: Close Cabinet Bottom Drawer	0.1574	0.0741	1.9118	0.6659	0.0057	4.59%
Task 9: Place the black bowl into the cabinet’s bottom drawer	0.1551	0.0913	1.8665	1.3193	0.0161	2.03%
Task 10: Place Wine Bottle on Wine Rack	0.1841	0.1114	2.0298	1.0022	0.0625	4.86%
Task 11: Close Cabinet Top Drawer	0.1224	0.0683	1.5424	0.5822	0.0053	2.62%
Task 12: Place the black bowl into the cabinet’s top drawer	0.2109	0.1418	2.6997	2.1633	0.0961	11.4%
Task 13: Place Black Bowl on Plate	0.1682	0.0816	2.2262	1.1552	0.2511	2.14%
Task 14: Place the black bowl on top of the cabinet	0.1375	0.0736	1.8499	0.9924	0.0108	6.73%
Task 15: Place the right moka pot on the stove	0.1890	0.1130	2.5768	1.4929	0.0134	3.19%
Task 16: Turn off the stove	0.2149	0.1213	2.8268	1.5823	0.0219	8.93%
Average	0.2128	0.0985	2.7844	1.4670	0.0398	4.79%

681 object states required for timestep-by-timestep replay. This distinction matters because an action
682 rollout can diverge from the originally verified trajectory under small simulator, contact, or controller
683 differences, while state replay keeps the rendered visual stream aligned with the successful task
684 execution.

685 Task admission is conservative. A task must have enough verified attempts both remotely and lo-
686 cally, and each selected attempt must be long enough to support action-chunk training. Candidate
687 attempts are ranked by shortest verified simulation time and then by attempt id, which biases the
688 rendered corpus toward compact successful demonstrations. The default setting renders all eligible
689 attempts for an admitted task, so the data-generation stage emphasizes coverage once the verification
690 threshold has been met.

691 Every rendered episode contains two RGB observation streams: a third-view camera that observes
692 the workspace from outside the robot and a wrist camera mounted near the end-effector. This two-
693 view design is used consistently throughout rendering, AXIS pretraining, and LIBERO post-training.

694 During rendering, the simulator state is set explicitly at each timestep. The replay record includes
695 robot joint position, joint velocity, joint position target, object position, object rotation, object lin-
696 ear velocity, and object angular velocity. The joint position target is retained because it supports
697 downstream action supervision and consistency checks, while the object state fields are required to
698 reconstruct task-relevant interactions without relying on a fresh dynamics rollout.

699 In the IsaacSim rendering path, physics stepping is disabled during replay. The replay state is
700 therefore authoritative: the renderer refreshes cameras and sensors after the state is set, but it does
701 not use a new physics integration step to determine the next state. This makes the data-generation
702 process closer to visualizing a verified trajectory than to resimulating a controller. It also explains
703 why table-height or scene-geometry randomization must be coupled with corresponding adjustments

Table 8: Default rendering, camera, and domain-randomization parameters.

Category	Parameter	Value
Rendering backend	Simulator	IsaacSim / IsaacLab
Rendering backend	Render mode	Ray-traced lighting
Rendering backend	Render resolution	256 × 256
Rendering backend	Camera data type	RGB
Rendering backend	Depth	Disabled by default
Robot	Robot model	Franka, fixed base, position control
Third-view camera	Position	(1.21, 0.0, 0.885)
Third-view camera	Look-at target	(0.3, 0.0, 0.0)
Third-view camera	Focal length	24.0
Third-view camera	RGB field of view	70° horizontal, 43° vertical
Third-view camera	Focus distance	400.0
Third-view camera	Clipping range	(0.05, 100000.0)
Wrist camera	Mount link	Franka hand
Wrist camera	Local position	(0.16, 0.0, 0.01)
Wrist camera	Local quaternion	(0.0, -0.38268343, 0.0, -0.92387953)
Wrist camera	RGB field of view	87° horizontal, 58° vertical
Third-view camera randomization	Position delta	±0.10 m along each axis
Third-view camera randomization	Roll/pitch/yaw delta	±8°
Wrist camera randomization	Depth / lateral / vertical delta	±0.02 m, ±0.015 m, ±0.005 m
Wrist camera randomization	Roll / pitch / yaw delta	±8°, ±4°, ±4°
Scene randomization	Randomization level	Level 3
Scene randomization	Scene mode	Mode 3, full USD scene
Lighting randomization	Main light intensity	12,000–35,000
Lighting randomization	Corner light intensity	5,000–15,000
Lighting randomization	Main color temperature	2,800–6,500 K
Lighting randomization	Sphere color temperature	2,500–6,000 K
Lighting randomization	Disk orientation delta	±15°
Lighting randomization	Sphere position delta	±(0.5, 0.5, 0.3) m

704 to replayed robot/object states and camera look-at targets. The randomized scene changes the visual
705 and geometric context, but the replayed trajectory must remain physically placed on the randomized
706 table.

707 Validation is performed at the demonstration level. A demonstration is rejected if required RGB
708 videos or metadata are missing, if joint traces are empty, non-finite, all zero, length-mismatched,
709 or too short for action-chunk training. The validated output contains the third-view RGB stream,
710 wrist RGB stream, metadata, joint-position trace, joint-position-target trace, and episode boundary
711 information. These elements are then packaged into training-ready records with third-view RGB,
712 wrist RGB, state, and 7D LIBERO-style joint-position actions.

713 The default rendering backend uses IsaacSim / IsaacLab with ray-traced lighting. Depth is disabled
714 by default, so the stable visual data consist of third-view RGB and wrist RGB. The third-view
715 camera approximates a fixed external RGB sensor, while the wrist camera is mounted on the Franka
716 hand and approximates an eye-in-hand sensor. Both camera pose randomization and full scene
717 randomization are enabled.

718 The randomization settings are broader than simple camera jitter. Scene mode 3 enables full USD
719 scene composition, and level 3 includes scene, material, lighting, and camera variation. The intended
720 effect is to increase visual diversity while keeping the underlying state/action supervision tied to
721 verified trajectories. Because depth is disabled by default, most of the rendered data volume comes
722 from RGB videos, metadata, and wrist RGB observations.

723 **F** $\pi_{0.5}$ **Pretraining from Base**

724 AXIS pretraining uses verified episodes converted into a LIBERO-style visual imitation-learning
 725 format. The default snapshot follows the paper-wide AXIS release scale: 207 tasks and 50,129
 726 episodes. Each episode pairs continuous robot state/action supervision with a third-view camera, a
 727 wrist camera, and episode-boundary information.

Table 9: AXIS-Datasets snapshot and default $\pi_{0.5}$ input dimensions.

Dataset / Model Parameter	Value
AXIS task count	207
AXIS episode count	50,129
Observation cameras	Third-view RGB and wrist RGB
Image resolution	256 × 256
State dimension	8D robot state
Action dimension	7D joint-position action
Action horizon	10
Model action dimension after padding	32
Model state dimension after padding	32

728 Training samples are formed as observation-action chunks within individual episodes. For each
 729 selected timestep, the sampler reads the current state and image observations, then constructs a 10-
 730 step action chunk from the same episode. If the timestep is close to an episode boundary, the action
 731 targets are clipped at the episode end. This prevents action supervision from crossing into the next
 732 demonstration and makes the action target semantically consistent with the current observation.

733 The training representation is intentionally simple: each sample contains robot state, a 10-step action
 734 target, third-view RGB, wrist RGB, and the episode boundary needed to keep action chunks within a
 735 single demonstration. This representation is then mapped into the LIBERO-compatible model input
 736 format.

737 Pretraining initializes from the official $\pi_{0.5}$ base JAX parameters and performs full-model finetun-
 738 ing. The model uses a PaliGemma Gemma-2B backbone and a Gemma-300M action expert. No
 739 LoRA adapter is used: the vision encoder, language backbone, and action expert are all trainable.
 740 $\pi_{0.5}$ mode is enabled, the action horizon is 10, and the tokenizer uses a maximum prompt length of
 741 200 tokens.

742 The input view set is intentionally narrow and matches the rendered dataset: the third-view image
 743 provides the external scene view, and the wrist image provides eye-in-hand context. State remains a
 744 continuous model input rather than being discretized into the language prompt. This keeps language
 745 reserved for task instructions while preserving low-level robot state in the continuous observation
 746 pathway.

747 The transform stack first maps third-view and wrist observations into the LIBERO-style observation
 748 layout, standardizes image layout, and applies quantile normalization to state and actions. Images
 749 are pad-resized from 256 × 256 to 224 × 224. During model loss computation, visual augmentation
 750 further applies random crop to 95%, rotation in $[-5^\circ, 5^\circ]$, and brightness/contrast/saturation jitter.
 751 Prompts are tokenized with the PaliGemma SentencePiece tokenizer. Finally, state is padded from
 752 8D to 32D and action chunks are padded from 7D to 32D; inference slices the output back to the
 753 first seven action dimensions.

754 For $\pi_{0.5}$, state and action normalization use quantile statistics:

$$x_{\text{norm}} = \frac{x - q_{01}}{q_{99} - q_{01} + 10^{-6}} \times 2.0 - 1.0.$$

755 The statistics are computed over state and actions. Action statistics use the same 10-step horizon as
 756 training, including clipping at episode boundaries. Images are not normalized through these statis-

757 tics; they are converted from uint8 image values into the model’s image range during observation
 758 construction.

759 The pretraining loss is flow matching over action chunks, not one-step behavior cloning. Given
 760 normalized and padded actions with shape [batch, 10, 32], the training step samples Gaussian noise
 761 and a continuous timestep, forms a noisy action suffix, and trains the model to predict the velocity
 762 target:

$$763 \quad \epsilon \sim \mathcal{N}(0, 1), \quad t \sim \text{Beta}(1.5, 1) \times 0.999 + 0.001,$$

$$x_t = t\epsilon + (1 - t)a, \quad u_t = \epsilon - a.$$

764 The model receives observation prefix tokens and noisy action suffix tokens, then predicts the action
 765 velocity for the final action-horizon positions. The loss is mean squared error between the predicted
 766 velocity and u_t , averaged over the action horizon and action dimensions. Since actions are padded to
 767 32 dimensions, the unused dimensions become fixed zero targets unless explicit action-loss weights
 768 are introduced. Thus, the model is optimized over the full padded action interface even though only
 769 the first seven dimensions are returned at inference time.

Table 10: Default $\pi_{0.5}$ AXIS pretraining optimization settings.

Training Parameter	Value
Global batch size	8
Training hardware	8×A100 GPUs
Train steps	100,000
Optimizer	AdamW
AdamW β_1	0.9
AdamW β_2	0.95
AdamW ϵ	10^{-8}
Weight decay	10^{-10}
Gradient clipping	1.0 global norm
Peak learning rate	5×10^{-5}
Warmup steps	10,000
Post-warmup learning rate	Effectively constant at 5×10^{-5}
EMA decay	0.999

770 Training length is reported in optimization steps rather than epochs, because each sample is
 771 an observation-action chunk drawn from within an episode. Checkpoints contain training state,
 772 inference-facing parameters, and dataset normalization artifacts. When EMA is enabled, the
 773 inference-facing parameters use EMA-smoothed weights and are the intended source for down-
 774 stream LIBERO post-training initialization.

775 **G LIBERO Post-training from AXIS-Pretrained Parameters**

776 LIBERO post-training starts from AXIS-pretrained $\pi_{0.5}$ inference parameters. Only model paramete-
 777 rs are transferred. Optimizer state, training step, and EMA state are initialized for a new LIBERO
 778 run. This is different from resuming training: resume is only appropriate when continuing the same
 779 LIBERO run, because it restores the full training state rather than using the AXIS checkpoint as an
 780 initialization source.

781 The AXIS-to-LIBERO transfer also keeps dataset normalization separate. The AXIS checkpoint
 782 provides weights, but LIBERO post-training uses LIBERO-specific normalization statistics stored
 783 with the post-training checkpoint. This prevents the LIBERO policy from using AXIS state/action
 784 normalization at serving time. The AXIS-to-LIBERO setup uses the final AXIS full-pretraining
 785 parameters as initialization and trains the LIBERO stage for 30,000 steps.

786 LIBERO post-training uses the official LIBERO LeRobot dataset. The data provide third-view cam-
 787 era images, wrist camera images, 8D continuous policy state, task language prompts, and native
 788 LIBERO delta actions. Because LIBERO actions are already delta actions, no additional delta con-

789 version is applied. This keeps the post-training visual input consistent with the AXIS-pretraining
 790 setup.

791 The post-training input representation remains aligned with $\pi_{0.5}$. The third-view camera and wrist
 792 camera provide visual observations, state remains an 8D continuous vector before padding, and the
 793 action target is a 10-step chunk. The output transform returns the first seven action dimensions.
 794 LIBERO-specific normalization statistics are stored with the post-training checkpoint and are re-
 795 quired for downstream policy loading.

Table 11: Default LIBERO post-training data, input, and optimization settings.

LIBERO / Post-training Parameter	Value
Dataset	Official LIBERO LeRobot release
Expected episodes	1,693
Image shape	[256, 256, 3]
State dimension	8D LIBERO policy state
Action horizon	10
Action semantics	Native LIBERO delta action
Extra delta conversion	Disabled
Active cameras	Third-view camera and wrist camera
Train steps	30,000
Global batch size	64
Training and evaluation hardware	8×A100 GPUs
Optimizer	AdamW
Gradient clipping	1.0 global norm
Peak learning rate	5×10^{-5}
Warmup steps	10,000
Post-warmup learning rate	Effectively constant at 5×10^{-5}
EMA decay	0.999

796 The default 8×A100 setting uses global batch size 64, which preserves approximately eight sam-
 797 ples per GPU. The optimization setup mirrors the pretraining learning-rate and EMA choices while
 798 changing the data distribution and initialization.

799 A LIBERO post-training checkpoint must include inference parameters and LIBERO-specific nor-
 800 malization artifacts. Loading only the model weights is insufficient for downstream policy con-
 801 struction because normalization statistics are required to build the trained policy correctly. The
 802 checkpoint therefore stores inference-facing model weights, training state for resuming the same
 803 run, and dataset-specific normalization statistics.

804 H Real-time Inference and Real-world Rollouts

805 This section supplements the main paper’s simulation-based LIBERO-Plus evaluations with qualita-
 806 tive real-world rollouts. We evaluate real-time inference on a Franka robot equipped with a gripper.
 807 For the qualitative rollouts, we selected two tasks: grasping and pick-and-place.

808 During real-time inference, the policy receives RGB observations from a global camera and a wrist-
 809 mounted camera, together with the continuous robot state. The policy predicts a short-horizon action
 810 chunk, which is converted into executable commands for the Franka controller at deployment time.
 811 This setup uses the same types of visual and proprioceptive inputs as the training pipeline, while
 812 leaving calibration, timing, and low-level controller execution to the real-robot system.

813 In addition to real-world keyframes, we also include a qualitative comparison between simulation
 814 replay and real-world execution. The simulation row visualizes the replayed task trajectory, while
 815 the real-world rows show the corresponding view observations. This comparison is intended to
 816 illustrate the observation-level correspondence between the replay pipeline and the real-robot setup.
 817 These examples are intended as qualitative evidence of data usability and policy executability.

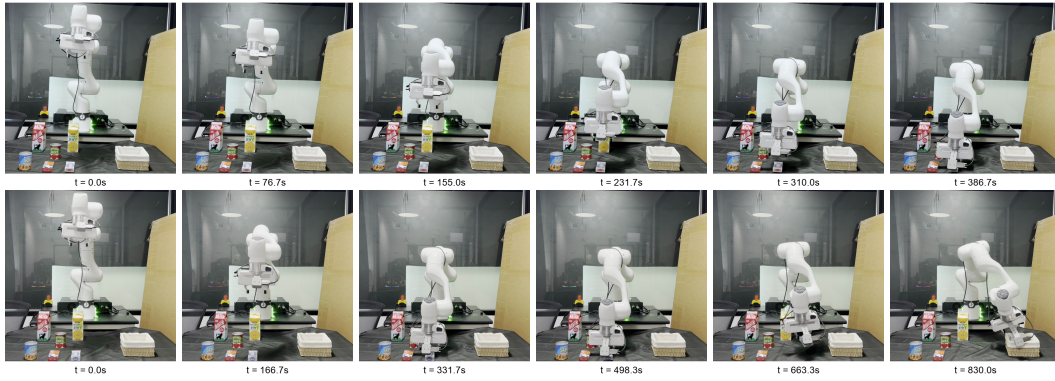


Figure 5: Real-world Franka rollouts. The first row shows a grasping task, and the second row shows a pick-and-place task. Timestamps indicate elapsed real execution time.



Figure 6: Real-time pick-and-place rollout. The two rows show the global and wrist-view.

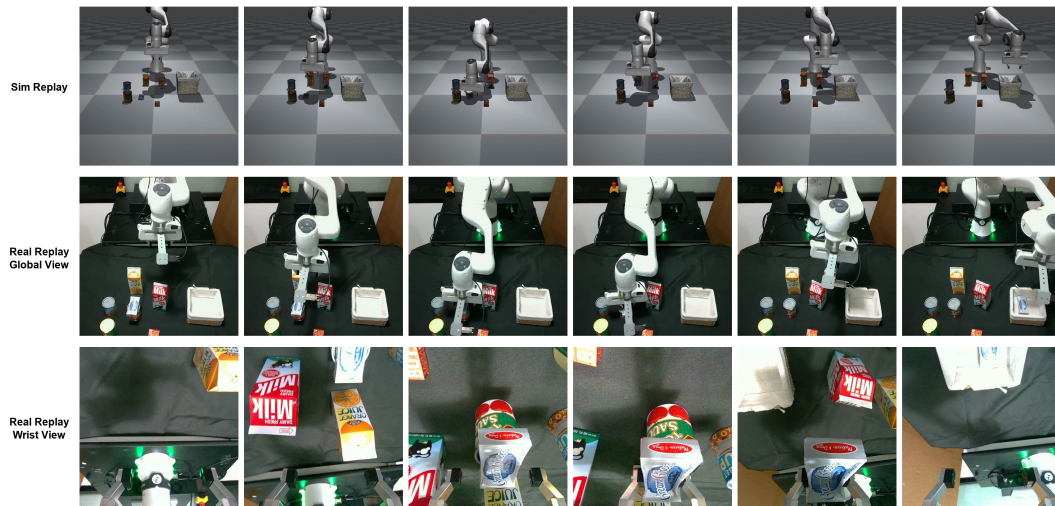


Figure 7: Pick-and-place replay. The three rows show the simulation replay, real-world global view, and real-world wrist view, respectively.