
AnnotateAnything: Automatic Annotation of 3D Assets for Robot Manipulation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Simulation enables scalable robot data collection, but raw 3D assets provide only ge-
2 ometry and appearance, lacking the semantic, interactive, and physical knowledge
3 needed to specify where and how robots should act. In this work, we present **Anno-**
4 **teAnything**, a general automatic annotation framework that converts passive 3D
5 assets into manipulation-ready assets with structured, diverse, and executable ma-
6 nipulation labels. AnnotateAnything is built around two complementary pipelines.
7 First, a unified **visual-language annotation pipeline** leverages vision-language
8 reasoning to infer object semantics, interaction constraints, and 3D-grounded cues
9 such as keypoints, and parts, providing human-prior guidance for identifying mean-
10 ingful interaction regions. Second, a fully automatic and massively parallel **physics**
11 **annotation pipeline** grounds these priors in each asset’s geometry, kinematics,
12 and physical constraints through candidate generation, geometry optimization,
13 trajectory generation and data augmentation. This pipeline produces diverse and ex-
14 ecutable action annotations, including grasp poses, dexterous contacts, articulation
15 waypoints, insertion directions, hanging affordances, and navigation targets. Based
16 on the generated annotations, we further build an asynchronous parallel simulation
17 data-collection system across diverse objects, tasks, and robot embodiments. Exper-
18 iments demonstrate that AnnotateAnything achieves superior annotation efficiency,
19 data-collection efficiency, and task success rates over existing annotation and data-
20 generation pipelines, while also supporting downstream tasks such as affordance
21 detection, robotic VQA, and visual instruction finetuning. We will open-source our
22 code, annotations, and benchmark to facilitate future research. Project page: XXXX.

23 1 Introduction

24 Training general-purpose robotic agents requires large amounts of diverse interaction data [1–5].
25 Recent progress in imitation learning and robot foundation models has further highlighted the
26 importance of scaling robot data across objects, scenes, and tasks [6–12]. However, collecting real-
27 world robot trajectories remains expensive, labor-intensive, and difficult to scale [13–16]. Simulation
28 provides a promising alternative, where interactions can be generated automatically, safely, and in
29 parallel [17–20]. Recent asset repositories and simulation benchmarks further provide rich geometric
30 and visual substrates for simulation-based data generation [21–24]. Yet assets and environments alone
31 are not enough: raw 3D assets usually describe only shape and appearance, but lack the semantic,
32 interactive, and physical knowledge that makes robot data valuable, such as where to grasp, which
33 part can move, how an object should be articulated, or where insertion and hanging are feasible.
34 Therefore, the central bottleneck is not merely obtaining more assets, but automatically producing
35 manipulation annotations that turn passive 3D geometry into actionable robot interaction data.



Figure 1: Teaser. High-level overview of AnnotateAnything and its visual-language-action annotation pipeline.

36 Generating manipulation annotations for raw 3D assets presents four key challenges. First, useful
 37 annotations must encode human interaction priors so that the resulting robot trajectories are human-
 38 like, safe, and likely to succeed. *This requires automatic annotation to incorporate high-level*
 39 *human-prior reasoning and interaction understanding (C1)* [25, 26]. Second, semantic annotation
 40 alone is insufficient: labels must be grounded in each asset’s geometry and kinematics, since feasible
 41 actions depend on shape, collision constraints, and articulated structure. *The challenge is to transform*
 42 *abstract interaction intent into asset-specific executable labels (C2)* [27, 28]. Third, diversity is
 43 essential for task deployment, as manipulation often admits multiple valid solutions across grasps,
 44 poses, and interaction strategies. *Automatic annotation must thus preserve diverse feasible labels*
 45 *rather than collapse to a single canonical solution (C3)* [29, 30]. Finally, annotations must scale
 46 across large numbers of assets, categories, and tasks while preserving each asset’s geometric precision
 47 and interaction specificity. *This calls for a fully scalable and automated process that avoids per-asset*
 48 *manual design or task-specific engineering (C4).*

49 Producing annotations that meet these requirements remains difficult. Existing works often rely on
 50 human annotation to provide semantic, interactive, or physical knowledge [31, 30, 32, 33], which is
 51 costly, labor-intensive, hard to scale, and limited in diversity. RL-based alternatives reduce manual
 52 labeling but require careful reward engineering, complex training, and substantial computation, while
 53 often producing behaviors misaligned with human priors and hard to deploy in realworld [34, 35].

54 Therefore, we present **AnnotateAnything**, a general automatic annotation framework that transforms
 55 raw 3D assets into manipulation-ready assets with structured, actionable, and executable robot ma-
 56 nipulation labels. AnnotateAnything provides three complementary types of annotations: language
 57 annotations for semantic and interaction-level reasoning, visual annotations for 3D-grounded under-
 58 standing, and action annotations that specify executable manipulation parameters such as grasp poses,
 59 articulation waypoints, and navigation trajectory.

60 To inject human interaction priors, AnnotateAnything builds a vision-language annotation pipeline
 61 that leverages VLMs to bridge language reasoning and visual grounding. Language annotations
 62 provide multi-granularity reasoning from object semantics to interaction constraints, while visual
 63 annotations translate such reasoning into 3D-grounded cues such as affordances, keypoints, and
 64 part segmentation. These cues localize meaningful interaction regions and guide the generation of
 65 action annotations, incorporating VLM-derived human priors and interaction understanding into
 66 automatic annotation (**tackling C1**). To ground annotations in each specific asset, AnnotateAnything
 67 combines geometry-based optimization with asset-specific physical reasoning. Geometry optimization
 68 refines grasps, contacts, and interaction regions under local shape and collision constraints, while
 69 physical reasoning derives waypoints and motion directions from the asset’s kinematic structure
 70 and physical properties. This converts abstract priors into precise, asset-specific executable labels
 71 (**tackling C2**)[36–38]. To preserve diversity, AnnotateAnything generates annotations at three levels.
 72 Within each task, it produces many feasible candidates with different contacts, approach directions,

Table 1: Comparison of automatic data collection and annotation generation methods. ✓: explicitly supported or reported; ○: partially supported; ✗: not supported or not reported.

| Method | Skill / Interaction Coverage | | | | | | | Scenario / Embodiment | | | Pipeline Property | | |
|-------------------------|------------------------------|----------|---------|------------|--------------|------------|------------|-----------------------|--------|----------|--------------------|---------------------|-------------------|
| | Grasp | DexGrasp | BiGrasp | BiDexGrasp | Articulation | Deformable | Navigation | Tabletop | Mobile | Humanoid | Physics Validation | Parallel Generation | Data Augmentation |
| InternData-A1 [33] | ✓ | ✗ | ✓ | ✗ | ○ | ○ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| InternData-N1 [44] | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ○ | ✗ | ✓ |
| RoboTwin2.0 [31] | ✓ | ✗ | ✓ | ✗ | ○ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ |
| MolmoBot [45] | ✓ | ✗ | ✗ | ✗ | ○ | ✗ | ○ | ✓ | ○ | ✗ | ✓ | ✗ | ✓ |
| RoboGen [34] | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ○ | ✗ | ○ | ✓ | ✗ |
| GenManip [46] | ✓ | ✗ | ✗ | ✗ | ○ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| OmniManip [47] | ✓ | ✗ | ✗ | ✗ | ○ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| GenieSim 3.0 [48] | ✓ | ○ | ○ | ○ | ○ | ✗ | ✗ | ✓ | ○ | ✗ | ✗ | ✗ | ✗ |
| AnnotateAnything | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

orientations, and interaction configurations, rather than a single canonical label. Across tasks, it covers primitives such as grasping, dexterous manipulation, bimanual interaction, articulation, folding, hanging, and mobile manipulation. Across embodiments, it supports parallel-jaw grippers, dexterous hands, dual-arm systems, mobile manipulators, and humanoids (**tackling C3**) [39, 40]. To enable scalability, AnnotateAnything builds a complete and fully parallelized physics annotation pipeline, covering candidate generation, geometry optimization, label selection, trajectory generation, physical validation, and data augmentation without per-asset manual design. With CUDA-accelerated optimization and validation [36], the pipeline scales across assets and tasks while preserving asset-specific accuracy, interaction feasibility, and high execution success rates (**tackling C4**).

Based on the generated annotations, we set up an asynchronous parallel data-collection system for efficient simulation-based robot data collection. Beyond data collection, our integrated visual-language-action annotations further support several downstream applications, including affordance and keypoint detection [41–43], robotics-oriented VQA and reasoning, and visual instruction finetuning for multimodal robot models. Experiments show that AnnotateAnything achieves substantially higher annotation efficiency, data-collection efficiency, and task success rates than existing annotation or data-generation pipelines, demonstrating the practical value of automatic, diverse, and asset-specific manipulation annotations.

In summary, our contributions are as follows:

- We present **AnnotateAnything**, a general **automatic annotation framework** that converts passive 3D assets into **manipulation-ready assets** through structured, actionable, and executable annotations.
- We design a unified **visual-language-action annotation pipeline** that integrates language reasoning, 3D visual grounding, and asset-specific action label generation, enabling **human-prior reasoning**, **3D-grounded understanding**, and **executable manipulation specification** across diverse objects, tasks, and robot embodiments.
- We develop a fully automatic and massively parallel **physics annotation pipeline**, covering candidate generation, geometry optimization, label selection, trajectory generation, physical validation, and data augmentation, with CUDA-accelerated modules for scalable annotation without sacrificing **diversity**, **feasibility**, or **asset-specific precision**.
- We demonstrate **superior annotation efficiency**, data-collection efficiency, and task success rates over existing methods, together with utility across several downstream tasks.

2 Related Work

Automatic Data Collection for Robot Learning. Automatic data collection mainly follows two paradigms. Annotation-based methods generate structured, often human-like trajectories from demonstrations, scripts, or manually specified task annotations [33, 31, 49], but are limited by human effort and task-specific design. RL-based methods automate data generation in simulation [35, 34, 50–53], but require costly training, reward engineering, and careful reset or success definitions, while often producing behaviors that are not human-like. In contrast, AnnotateAnything automatically generates skill-consumable annotations from raw 3D assets, enabling reusable atomic skills to collect structured and scalable robot interaction data.

113 **Automatic Annotation for Manipulation** Existing automatic annotation methods often formulate
 114 manipulation annotation as task-specific optimization, such as generating grasps, contact points,
 115 affordances, or motion waypoints for a particular skill. While effective in specific settings, these
 116 methods usually cover limited object diversity and annotation scale[29, 46], suffer from low success
 117 rates on complex assets[50, 51, 38, 54], or support only a narrow set of embodiments and task
 118 types[55, 53]. AnnotateAnything provides a more general annotation pipeline that produces diverse
 119 skill-consumable labels across grasping, dexterous manipulation, bimanual interaction, articulation,
 120 insertion, and hanging, enabling scalable data collection over broad 3D assets and manipulation skills.
 121 Detailed discussion please refer to Appendix F

122 3 Method

123 AnnotateAnything takes as input a raw 3D asset, ranging from a single object to a room-scale
 124 scene, and converts it into a manipulation-ready asset with hierarchical language, visual, and action
 125 annotations. Our framework first leverages VLMs to infer human interaction priors and ground
 126 them into multi-level asset and scene annotations, including sparse/dense descriptions, keypoints,
 127 part segmentation, occupancy maps, floor plans, and top-view layouts (Sec. 3.1). These annotations
 128 are generated through hierarchical language reasoning (Sec. 3.1.1), hierarchical visual grounding
 129 (Sec. 3.1.2), and cross-level asset-to-room composition (Sec. 3.1.3). The physics-based pipeline then
 130 converts these visual-language cues into executable action annotations across rigid objects, articulated
 131 objects, garments and deformable objects, and room-scale scenes (Sec. 3.2), using a unified action
 132 schema (Sec. 3.2.1), a general generation-and-validation pipeline (Sec. 3.2.2), and primitive-specific
 133 instantiations (Sec. D).

Table 2: Hierarchical annotation schema of AnnotateAnything.

| Level | Annotation Type | Example Outputs |
|--------|--------------------|--|
| Asset | Language | Sparse tags; dense object-, part- and task-level descriptions |
| Asset | Visual | 3D keypoints; part segmentation; |
| Room | Language | Scene descriptions; object relations; task contexts |
| Room | Visual | Occupancy maps; floor plans; top-view/BEV maps; object layouts |
| Action | Rigid object | Parallel-jaw grasps; dexterous contacts; insertion and hanging poses |
| Action | Articulated object | Handles; motion axes; articulation waypoints; opening and closing trajectories |
| Action | Deformable object | Garment keypoints; bimanual grasp points; folding and hanging trajectories |
| Action | Room-scale scene | Navigation targets; approach poses; interaction-ready candidate base poses |

134 3.1 Hierarchical Visual-Language Annotation Pipeline

135 3.1.1 Hierarchical Language Annotation

136 **Asset-level language annotation.** For each object-level asset, we render multi-view RGB observa-
 137 tions and use Qwen3-VL [56] to generate **three-level language annotations**: a semantic phrase, a
 138 functional sentence, and a dense paragraph for **part-aware interaction reasoning**. The dense anno-
 139 tation grounds segmented parts with their functions and feasible interactions, guiding downstream
 140 visual grounding and action generation.

141 **Room-level language annotation.** Given selected room views and the labeled top-view occupancy
 142 map, we use Qwen3.5-VL [56] to produce **three-level room annotations**: a layout summary, a
 143 furniture-and-zone description, and a dense paragraph covering object relations, navigable regions,
 144 interaction areas, and possible long-horizon robot tasks. These annotations complement asset-level
 145 labels with functional scene context.

146 3.1.2 Hierarchical Visual Annotation

147 **Asset-level visual annotation.** For each object-level asset, we reconstruct a fused 3D point cloud
 148 from multi-view RGB-D observations, making the pipeline consistent with real-world visual inference.
 149 We annotate **semantic keypoints** by sampling candidates with FPS and using the VLM to select
 150 functionally important anchors for interaction reasoning. In parallel, we generate **part segmentation**
 151 **annotations** with a native 3D pipeline based on Hunyuan3D-style decomposition, P3-SAM [57],

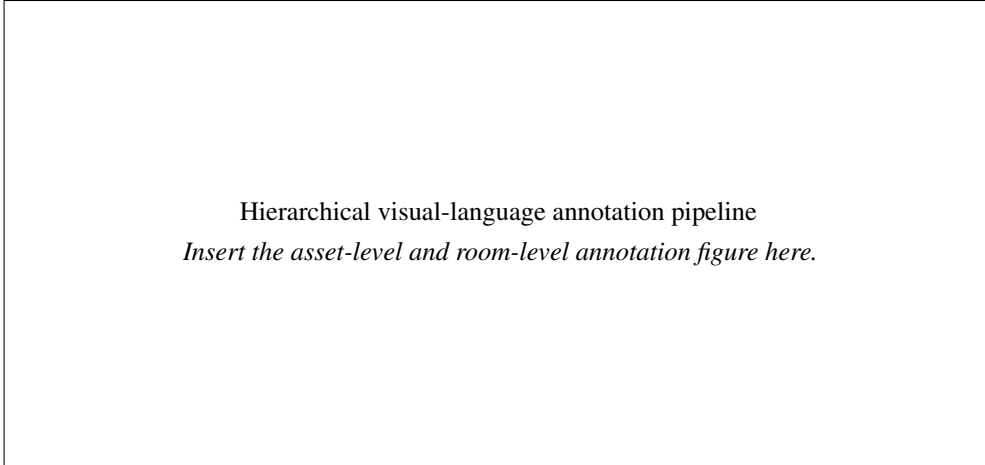


Figure 2: Hierarchical visual-language annotation pipeline. AnnotateAnything generates asset-level language and visual annotations, including sparse/dense descriptions, keypoints, part segmentation, and affordance regions, and composes them into room-level annotations such as scene descriptions, occupancy maps, floor plans, top-view maps, and object layouts.

152 and X-Part [58]. These visual annotations provide 3D-grounded cues for downstream physics-based
153 action annotation.

154 **Room-level visual annotation.** For each room-scale scene, we generate visual annotations by
155 actively scanning the environment with simulated LiDAR and ray casting, and then stitching the
156 observations from individual rooms into a global scene representation. To capture spatial structure at
157 different levels, we construct multi-height occupancy maps, where each height slice encodes free
158 space, obstacles, and object geometry visible at that elevation. We further derive floor-plan and
159 wall-structure annotations by using VLM-based filtering to remove movable objects and furniture
160 from the raw occupancy evidence, retaining only persistent structural elements such as walls, doors,
161 and room boundaries. These room-level visual annotations provide global geometric context for
162 spatial exploration, navigation target generation, and scene-level action annotation.

163 3.1.3 Cross-level Composition and Consistency

164 We build room-level annotations by composing object-level annotations through scene instance
165 identities and 6D object poses. For each room-scale scene, we extract individual object instances
166 and process each one with the asset-level annotation pipeline, producing object-centric language
167 descriptions, keypoints, and part-masks. These annotations are then transformed from the object
168 coordinate frame to the global room coordinate frame and associated with room-level occupancy
169 maps, floor plans, top-view layouts, and scene descriptions. The resulting cross-level representation
170 links each object instance in the room to its semantic, visual, and affordance annotations, while
171 maintaining shared identifiers across language annotations, visual maps, and physics-based action
172 labels.

173 3.2 Physics-based Action Annotation Pipeline

174 The visual-language pipeline provides semantically meaningful 3D interaction regions, which are
175 further grounded in robot embodiment and physical feasibility by our physics-based action annotation
176 pipeline. This pipeline converts grounded visual-language priors into structured, executable, and
177 validation-aware action annotations through candidate target generation, trajectory generation, trajec-
178 tory optimization, physics validation, and physics-aware augmentation. The resulting manipulation-
179 ready assets contain diverse action labels, including grasp poses, dexterous contacts, articulation
180 waypoints, insertion directions, hanging poses, deformable-object trajectories, and navigation targets.

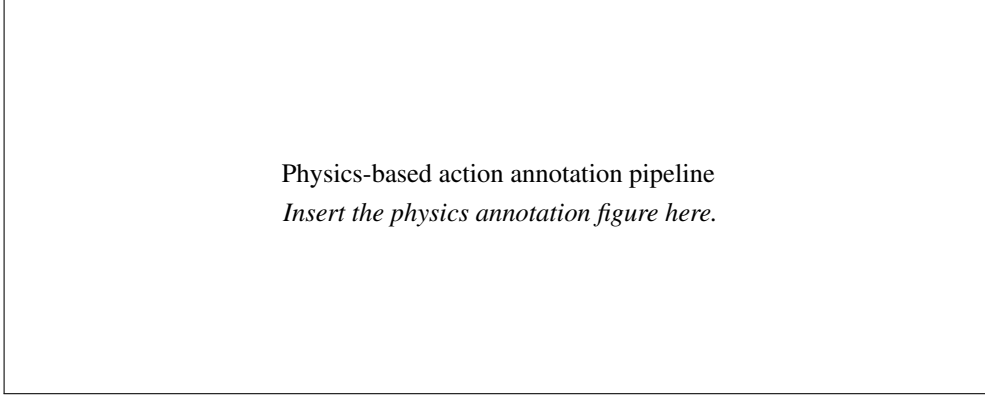


Figure 3: Physics-based action annotation pipeline. Grounded visual-language priors are converted into executable action annotations through candidate target generation, trajectory generation, trajectory optimization, physics validation, and physics-aware augmentation.

181 3.2.1 Unified Action Annotation Schema

182 A key design goal of our action annotation schema is to preserve **action diversity** while maintaining
 183 **functional consistency**. Instead of assigning a single canonical action to each object or part, we
 184 organize action annotations as a hierarchical candidate bank built on top of visual annotations. Given
 185 an asset \mathcal{A} , the visual annotation stage provides grounded interaction anchors

$$\mathcal{H}(\mathcal{A}) = \mathcal{P}(\mathcal{A}) \cup \mathcal{K}(\mathcal{A}) \cup \mathcal{R}^{\text{aff}}(\mathcal{A}) \cup \mathcal{R}^{\text{scene}}(\mathcal{A}),$$

186 where \mathcal{P} denotes part-level regions, \mathcal{K} denotes semantic keypoints, \mathcal{R}^{aff} denotes affordance regions,
 187 and $\mathcal{R}^{\text{scene}}$ denotes scene-level regions such as reachable free space or object-centric approach regions.
 188 Importantly, visual grounding provides not only the spatial anchor h , but also its functional affordance
 189 $\phi(h)$, such as grasp-for-lifting, grasp-for-opening, insertion, hanging, or navigation approach. Thus,
 190 parts and keypoints define local functional manipulation targets, while scene-level regions define
 191 navigation and approach targets.

192 For each anchor $h \in \mathcal{H}(\mathcal{A})$, we infer compatible skills $\mathcal{S}(h)$ from its functional affordance $\phi(h)$,
 193 language priors, geometry, and task constraints. This is crucial for long-horizon manipulation, where
 194 a physically stable grasp may be insufficient unless it supports the intended downstream interaction,
 195 e.g., grasping a mug handle for pouring or a cabinet handle for opening. We then maintain a **candidate**
 196 **bank** for every anchor-skill pair:

$$\mathcal{B}_{h,s} = \left\{ a_{h,s}^{(i)} \right\}_{i=1}^{N_{h,s}}, \quad s \in \mathcal{S}(h),$$

197 where $N_{h,s}$ is the number of feasible candidates for skill s at anchor h . Each candidate action is
 198 represented as

$$a_{h,s}^{(i)} = \left(s, o(h), h, \phi(h), x^{(i)}, \theta^{(i)}, \tau^{(i)}, v^{(i)}, d^{(i)} \right).$$

199 Here, s is the skill type, $o(h)$ is the associated object, part, or scene instance, and $\phi(h)$ conditions
 200 action generation with functional affordance. The concrete target $x^{(i)}$ is instantiated from the anchor,
 201 such as a contact point, surface patch, support edge, opening, articulated handle, or navigation base
 202 pose. The parameters $\theta^{(i)}$ store skill-specific quantities such as a 6D grasp pose, gripper width, contact
 203 set, hand configuration, insertion direction, hanging pose, or target base pose. The optional trajectory
 204 $\tau^{(i)}$ stores waypoints or motion templates for extended actions. The metadata $v^{(i)}$ records physical
 205 and functional feasibility, including collision status, IK feasibility, contact stability, task success,
 206 trajectory statistics, and affordance preservation. The descriptor $d^{(i)}$ records diversity factors such as
 207 target location, approach direction, contact mode, trajectory family, and embodiment parameters.

208 This schema defines action annotation as a **one-to-many mapping** from grounded visual anchors to
 209 function-conditioned executable candidates. Diversity exists across anchors and within each anchor-
 210 skill pair, since $\mathcal{B}_{h,s}$ stores multiple validated candidates with different poses, contacts, approach
 211 directions, trajectories, and feasibility scores. As a result, AnnotateAnything provides downstream

212 modules with a diverse bank of **physically feasible and functionally meaningful** action labels rather
213 than a single heuristic annotation.

214 3.2.2 General Physics-based Action Generation Pipeline

215 **Candidate target generation.** Given a skill type s , we first use language annotations and visual
216 grounding to localize skill-compatible candidate regions C_s from part masks, semantic keypoints,
217 affordance regions, or room-level occupancy/BEV maps. For example, functional grasping selects
218 grasp-for-use regions such as mug handles, articulation selects handles or movable parts, and nav-
219 igation selects traversable free-space regions around the target object. We then sample concrete
220 anchor targets from C_s using farthest-point sampling on object surfaces or grid-based sampling on
221 scene maps. Depending on the primitive, a sampled target can be a contact point, surface patch,
222 support edge, opening, articulated handle point, navigation goal, or interaction-ready base pose. The
223 sampled targets are filtered by geometric and embodiment constraints, including surface normals,
224 curvature, visibility, collision margins, reachability, and traversability, before being used to instantiate
225 the candidate banks.

226 **Trajectory generation.** Given a retained target $x^{(i)}$ for skill s , we generate an initial trajectory or
227 waypoint sequence

$$\tau_0^{(i)} = \mathcal{G}_s \left(h, \phi(h), x^{(i)}, \theta^{(i)}; \mathcal{A} \right),$$

228 where \mathcal{G}_s is selected by the skill’s constraint source. **Object-property-conditioned** generators use
229 static geometry or garment keypoints to produce grasp, dexterous grasp, fling, or fold templates.
230 **Object-trajectory-conditioned** generators follow articulated part motion, e.g., preserving the end-
231 effector–handle relation while interpolating a door or drawer joint. **Object-vector-conditioned**
232 generators align actions with insertion, hanging, pouring, or placement vectors. **Scene-trajectory-**
233 **conditioned** generators use occupancy or BEV maps to sample interaction-ready base poses and
234 compute coarse navigation paths with A^* , connecting free-space regions to manipulation targets.
235 These trajectories serve as template-level seeds for subsequent optimization and validation.

236 **Trajectory optimization.** The template trajectories are further refined by optimizing action param-
237 eters and waypoint sequences under contact, geometry, and embodiment constraints:

$$(\theta_*^{(i)}, \tau_*^{(i)}) = \arg \min_{\theta, \tau} E_{\text{task}} + E_{\text{contact}} + E_{\text{coll}} + E_{\text{kin}} + E_{\text{smooth}}.$$

238 Here, E_{task} encourages the trajectory to preserve the intended functional affordance, E_{contact} im-
239 proves contact alignment and stability, E_{coll} penalizes penetration and unsafe clearance, E_{kin} enforces
240 robot kinematics and embodiment constraints, and E_{smooth} regularizes waypoint motion. For dexter-
241 ous grasping, we optimize hand pose, finger configuration, and contact assignment so that the selected
242 contacts lie on the grounded functional region, satisfy joint limits and self-collision constraints,
243 and can support task-relevant object wrenches. For adaptive skills, we adjust waypoints according
244 to asset-specific geometry, e.g., scaling garment fold targets by sleeve length or hem-to-shoulder
245 distance. For navigation, the global path from A^* is locally optimized through DWB-style rollout
246 under embodiment-specific dynamics, such as differential-drive, Ackermann, or holonomic motion
247 constraints, producing feasible trajectories that respect clearance, turning radius, and manipulation
248 reachability.

249 **Physics validation.** Optimized candidates are validated through batched physics simulation before
250 being stored as action annotations. For most object-centric skills, we use floating embodiments, such
251 as a floating parallel gripper or floating dexterous hand, to validate contact, collision, stability, and
252 task progress independently of a specific robot arm or base pose. This decoupling avoids prematurely
253 rejecting valid object interactions, since downstream data collection may randomize robot-object
254 poses and embodiments. For strongly embodiment-dependent skills, such as garment pushing,
255 washing, or retrieval from entangled configurations, we instead validate with the full manipulator
256 to account for reachability, arm collision, and entanglement. To test robustness, we additionally
257 apply disturbances and vary gravity magnitude or direction during validation, especially for grasp and
258 dexterous grasp candidates. Accepted candidates must satisfy skill-specific success criteria and are
259 stored with metadata such as success flags, collision statistics, contact stability, task progress, rollout
260 length, disturbance setting, and validation embodiment. Because many generated candidates fail

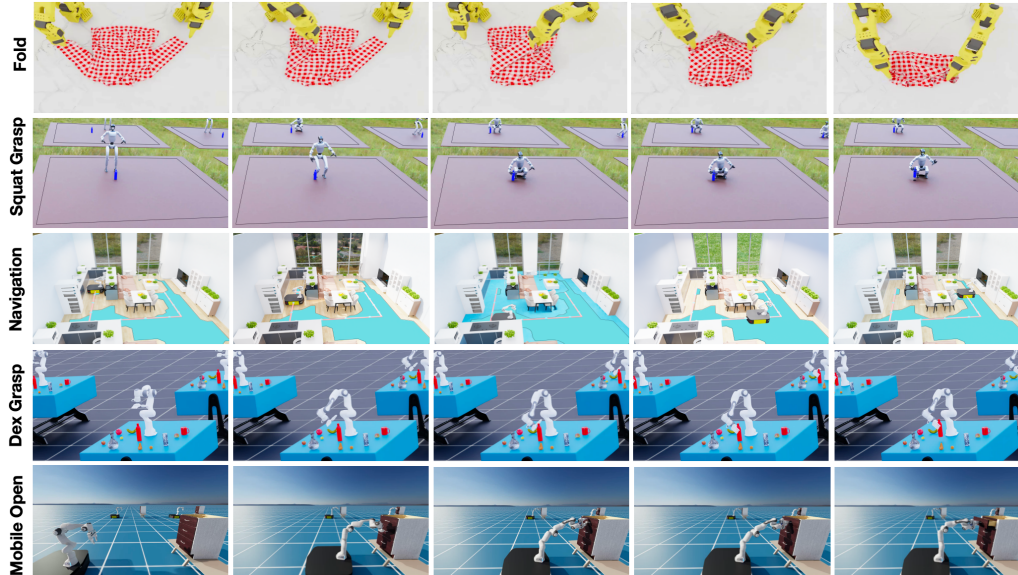


Figure 4: Atomic skill. Placeholder for the atomic-skill representation or taxonomy.

261 these checks, scalable parallel validation is critical for retaining high-quality labels and improving
 262 downstream success.

263 **Physics-aware augmentation.** To further increase diversity without breaking executability, we
 264 augment validated candidates in two ways. **Local perturbation** samples small variations around
 265 the validated anchor and action parameters while preserving the same functional affordance $\phi(h)$.
 266 For example, garment fold or fling points are jittered within local keypoint confidence regions,
 267 fold/place targets are perturbed within a bounded range, grasp candidates perturb contact centers,
 268 approach offsets, and in-plane rotations, and articulation candidates perturb handle contact points
 269 and pulling directions. These perturbations are bounded by part masks, affordance regions, collision
 270 margins, and embodiment constraints. **Symmetry-aware augmentation** uses symmetry priors from
 271 visual-language annotations to expand candidates on symmetric objects or parts. Let $\mathcal{G}_{\text{sym}}(o)$ denote
 272 the annotated symmetry group of object or part o . For a validated candidate $a_{h,s}^{(i)}$, we generate
 273 transformed candidates

$$\tilde{a}_{h,s}^{(i,g)} = g \cdot a_{h,s}^{(i)}, \quad g \in \mathcal{G}_{\text{sym}}(o),$$

274 by consistently transforming the target, pose parameters, and trajectory waypoints. For instance,
 275 grasp poses on a rotationally symmetric bottle can be rotated around its symmetry axis to produce
 276 additional feasible grasps. All augmented candidates are rechecked by lightweight geometric and
 277 physical constraints, and only feasible, affordance-preserving augmentations are retained in the
 278 candidate bank.

279 4 Downstream Tasks

280 AnnotateAnything turns passive 3D assets into reusable visual-language-action annotations for robot
 281 learning. As shown in Fig. 5, we study two downstream usages: using executable action annotations
 282 for large-scale simulation data collection, and converting the same annotations into supervision for
 283 affordance and keypoint detection, robot reasoning VQA, and 3D VLM instruction finetuning.

284 4.1 Large-scale Robot Data Collection

285 We use the generated action annotations as executable interfaces for large-scale simulation-based
 286 robot data collection. **This system is not the main contribution of this paper; we include it to**
 287 **demonstrate that AnnotateAnything annotations can be directly consumed by downstream**
 288 **robot execution and data-collection pipelines.**

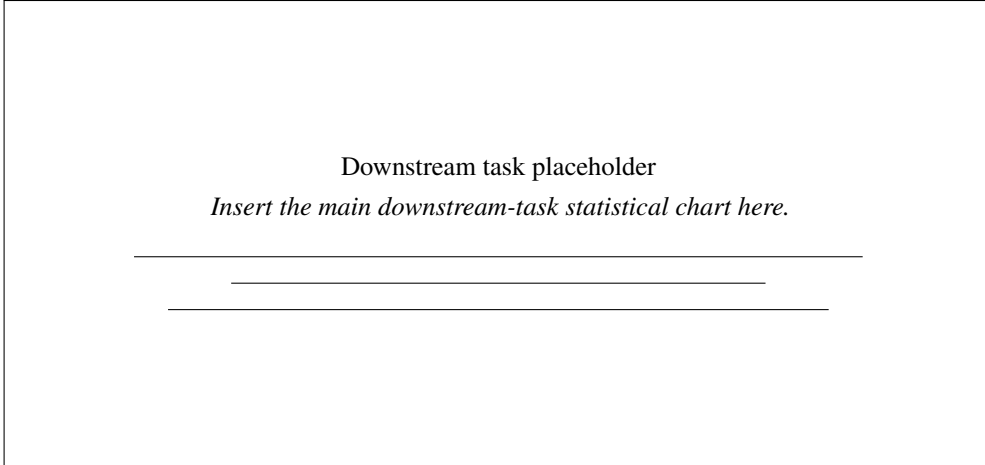


Figure 5: Downstream applications supported by AnnotateAnything. The same visual-language-action annotations can be converted into supervision for affordance and keypoint detection, robot reasoning VQA, and 3D VLM instruction finetuning.

289 **Atomic skill interface.** We build an atomic-skill library aligned with our annotation format, where
290 skills consume action parameters such as grasp poses, target parts, waypoints, insertion directions,
291 hanging anchors, and navigation goals. The library supports tabletop manipulation, bimanual
292 manipulation, whole-body control, humanoids, dexterous hands, and mobile manipulation, and can
293 compose atomic skills into long-horizon tasks.

294 **Parallel rollout generation.** We run heterogeneous parallel simulation environments, where each
295 environment independently samples assets, tasks, robot embodiments, object poses, and scene
296 layouts. For each rollout, asynchronous cuRobo-v2 planning is used for goal reaching and obstacle
297 avoidance [36, 37], with domain randomization over pose, lighting, material, camera, and scene
298 configuration.

299 **Trajectory validation.** For each randomized scene, we retrieve candidates from the annotation
300 bank, solve goal-set IK, and select the feasible solution with the lowest planning cost. We also remove
301 action candidates that repeatedly fail simulation validation, ensuring that the retained annotations
302 remain executable under diverse robot-object configurations.

303 4.2 Annotation-derived Downstream Applications

304 4.2.1 Keypoint Generation and Affordance Grounding

305 We reuse visual annotations as keypoint supervision and physics-validated action annotations as
306 affordance supervision. The former provides VLM-selected 3D functional anchors, while the latter
307 provides executable interaction labels, such as grasp directions, articulation trajectories, garment
308 pick-and-place points, and room-level navigation or interaction regions. These labels can be projected
309 to images, point clouds, or BEV maps for robot-centric perception training.

310 4.3 Robot Reasoning VQA

311 Robot reasoning VQA is constructed as a by-product of trajectory generation. During rollout, we
312 record execution-level ground truth, including the selected object, part anchor, skill type, IK goal-set
313 solution, approach side, navigation target, articulation direction, and garment manipulation points.
314 These labels depend on the runtime robot-object relative pose and the candidate selected from the
315 annotation bank, and therefore cannot be obtained from static assets alone. We use this mainly to
316 show that simulation rollouts generated from AnnotateAnything naturally provide grounded QA
317 supervision for robot execution reasoning.

Table 3: Visual-language annotation quality on the audited evaluation suite. Human scores are reported on a 0–100 scale. Detailed rubrics and sampling procedures are provided in the appendix.

| Method | Sem. | Ground. | Cover. | Action. | Overall |
|---------------|-------------|-------------|-------------|-------------|-------------|
| Direct VLM | XX.X | XX.X | XX.X | XX.X | XX.X |
| w/o 3D refine | XX.X | XX.X | XX.X | XX.X | XX.X |
| Ours | XX.X | XX.X | XX.X | XX.X | XX.X |

318 4.4 3D VLM Instruction-tuning Data

319 AnnotateAnything can also provide instruction-tuning data for 3D VLMs. In simulation, assets and
 320 rollouts provide dense grounding signals such as 3D bounding boxes, projected 2D boxes, instance
 321 masks, object poses, and part identities, which can be converted into grounding-style instruction-
 322 response pairs. Meanwhile, the language annotations and cross-level composition in Sec. 3.1.3
 323 provide object relations, room layouts, and task contexts, enabling QA pairs about 3D spatial
 324 relationships, object localization, and scene-level reasoning. We use this mainly as a lightweight
 325 downstream demonstration that our annotations can be repurposed into multimodal instruction data,
 326 rather than as a full 3D VLM benchmark.

327 5 Experiment

328 5.1 Experimental Setup

329 **Evaluation scope.** We evaluate AnnotateAnything as an automatic annotation pipeline rather than
 330 as a dataset release or a standalone data-generation system. Given heterogeneous raw 3D assets,
 331 our experiments test whether the pipeline can convert passive geometry into manipulation-ready
 332 language, visual, and action annotations that are semantically meaningful, 3D-grounded, physically
 333 valid, and useful for downstream simulation rollouts. We evaluate four aspects: annotation scale and
 334 conversion statistics, visual-language annotation quality, physics-grounded action annotation quality,
 335 and annotation-enabled data collection.

336 **Audited evaluation suite.** All quality and success-rate results are reported on an audited evaluation
 337 suite sampled from a larger heterogeneous asset pool. The pool spans multiple asset sources and
 338 asset types, including rigid objects, articulated objects, deformable or garment assets, and room-scale
 339 scenes. Since our goal is annotation conversion rather than curated asset release, we report source
 340 and scale statistics separately, while pass rates, success rates, and human scores are computed on the
 341 audited suite. The suite is stratified by valid asset–skill pairs, covering grasping, dexterous grasp-
 342 ing, bimanual grasping, bimanual dexterous grasping, articulation, insertion, hanging, deformable
 343 manipulation, and navigation or approach-target generation.

344 **Evaluation protocol.** Each experiment defines its metrics at the beginning of the corresponding
 345 subsection. Human evaluation is reported as an overall 0–100 score in the main paper, with detailed
 346 rubrics, sampling procedures, evaluator instructions, and per-dimension scores provided in the
 347 supplementary material.

348 5.2 Annotation Coverage over Heterogeneous Assets

349 5.3 Visual-Language Annotation Quality

350 We evaluate whether the visual-language stage produces annotations that are semantically correct,
 351 accurately grounded in 3D, and useful for subsequent action annotation. We compare **Direct**
 352 **VLM**, which directly predicts language and interaction regions from multi-view renderings; **w/o**
 353 **3D refinement**, which keeps hierarchical language reasoning but removes explicit 3D grounding
 354 refinement; and **Ours**, the full visual-language annotation pipeline. Human raters score sampled
 355 annotation bundles from the audited evaluation suite on a 0–100 scale. Table 3 reports semantic
 356 correctness, 3D grounding accuracy, coverage of interaction-relevant regions, actionability, and
 357 the overall score. Fig. ?? qualitatively shows how the predicted semantic priors are grounded into
 358 keypoints, parts, affordance regions, and scene-level spatial cues.

Table 4: Physics-grounded action annotation quality on the audited evaluation suite. The upper block reports results by action family, and the lower block reports macro-level ablations. Human scores are reported on a 0–100 scale. Full per-skill results, validation criteria, and human-evaluation rubrics are provided in the appendix.

| Setting | Phys. Pass | Asset Ready | Valid/Ready | Exec. Succ. | Ann./Hour | Human |
|------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| <i>Action-family results</i> | | | | | | |
| Grasping | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |
| Dexterous / Bimanual | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |
| Articulation | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |
| Insertion / Hanging | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |
| Deformable | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |
| Navigation / Approach | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |
| <i>Macro ablations</i> | | | | | | |
| Geometry-only | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |
| VL-only | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |
| w/o physics validation | – | XX.X | XX.X | XX.X | XX.X | XX.X |
| Ours full | XX.X | XX.X | XX.X | XX.X | XX.X | XX.X |

359 Ours achieves the highest overall score. Direct VLM prompting often produces plausible language
360 descriptions but is less reliable at localizing interaction regions in 3D. Removing 3D refinement
361 reduces grounding accuracy and actionability, showing that language reasoning alone is insufficient
362 for manipulation-ready annotation. These results indicate that the visual-language pipeline provides
363 reliable human-prior guidance for the physics annotation stage.

364 5.4 Physics-grounded Action Annotation Quality

365 We next evaluate whether grounded visual-language priors can be converted into executable action
366 annotations. For each valid asset–skill pair in the audited evaluation suite, the physics pipeline gener-
367 ates candidates, filters them by geometry and embodiment constraints, optimizes poses or trajectories,
368 and validates them in simulation. Table 4 summarizes the resulting annotation quality. The upper
369 block reports results by action family; the lower block reports macro-level ablations comparing
370 **Geometry-only**, **VL-only**, **w/o physics validation**, and **Ours**. We report physics validation pass rate,
371 asset readiness, valid annotations per ready asset–skill pair, execution success, annotation throughput,
372 and human score. Fig. ?? illustrates the final action labels produced by the physics pipeline.

373 The full pipeline produces validated annotations across all action families while maintaining high asset
374 readiness and execution success. The ablations show that geometry-only sampling lacks semantic
375 precision, VL-only annotation lacks physical executability, and removing physics validation reduces
376 downstream execution reliability. These results demonstrate that AnnotateAnything converts human-
377 prior visual-language cues into diverse, asset-specific, and physically executable action annotation
378 banks.

379 6 Conclusion

380 Do not change any aspects of the formatting parameters in the style files. In particular, do not modify
381 the width or length of the rectangle the text should fit into, and do not change font sizes. Please note
382 that pages should be numbered.

383 Most of the margin problems come from figures positioned by hand using `\special` or other
384 commands. We suggest using the command `\includegraphics` from the `graphicx` package.
385 Always specify the figure width as a multiple of the line width.

386 A number of width problems arise when L^AT_EX cannot properly hyphenate a line. Please give LaTeX
387 hyphenation hints using the `\-` command when necessary.

References

- 388
- 389 [1] Physical Intelligence. Pi-0.7: A steerable generalist robotic foundation model with emergent
390 capabilities. arXiv preprint, 2026. CorpusID: 287607456.
- 391 [2] NVIDIA. Gr00t n1: An open foundation model for generalist humanoid robots.
392 arXiv:2503.14734, 2025.
- 393 [3] Generalist AI Team. Gen-0: Embodied foundation models that scale with physical interaction.
394 Generalist AI Blog, 2025. November 4, 2025.
- 395 [4] Seonghyeon Ye, Yunhao Ge, Kaiyuan Zheng, and Joel Jang. World action models are zero-shot
396 policies. arXiv:2602.15922, 2026.
- 397 [5] Ruijie Zheng, Dantong Niu, Yuqi Xie, and Linxi Jim Fan. Egoscale: Scaling dexterous
398 manipulation with diverse egocentric human data. arXiv:2602.16710, 2026.
- 399 [6] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny
400 Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya
401 Ghosh, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming
402 Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl
403 Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Sta-
404 chowicz, James Tanner, Quan Vuong, Homer Rich Walke, Anna Walling, Haohuan Wang, Lili
405 Yu, and Ury Zhilinsky. π 0.5: a vision-language-action model with open-world generaliza-
406 tion. *ArXiv*, abs/2504.16054, 2025. URL [https://api.semanticscholar.org/CorpusID:
407 277993634](https://api.semanticscholar.org/CorpusID:277993634).
- 408 [7] Yue Chen, Muqing Jiang, Kaifeng Zheng, Jiaqi Liang, Chenrui Tie, Haoran Lu, Ruihai Wu,
409 and Hao Dong. Learning part-aware dense 3d feature field for generalizable articulated object
410 manipulation, 2026. URL <https://arxiv.org/abs/2602.14193>.
- 411 [8] Ruihai Wu, Haozhe Chen, Mingtong Zhang, Haoran Lu, Yitong Li, and Yunzhu Li. Neural
412 dynamics augmented diffusion policy. In *2025 IEEE International Conference on Robotics and
413 Automation (ICRA)*, pages 13234–13241, 2025. doi: 10.1109/ICRA55743.2025.11128651.
- 414 [9] Yan Shen, Ruihai Wu, Yubin Ke, Xinyuan Song, Zeyi Li, Xiaoqi Li, Hongwei Fan, Haoran Lu,
415 and Hao Dong. Biassemble: Learning collaborative affordance for bimanual geometric assem-
416 bly. *ArXiv*, abs/2506.06221, 2025. URL [https://api.semanticscholar.org/CorpusID:
417 279244917](https://api.semanticscholar.org/CorpusID:279244917).
- 418 [10] Yitong Li, Ruihai Wu, Haoran Lu, Chuanruo Ning, Yan Shen, Guanqi Zhan, and Hao Dong.
419 Broadcasting support relations recursively from local dynamics for object retrieval in clut-
420 ters. *ArXiv*, abs/2406.02283, 2024. URL [https://api.semanticscholar.org/CorpusID:
421 270226492](https://api.semanticscholar.org/CorpusID:270226492).
- 422 [11] Mingleyang Li, Yuran Wang, Yue Chen, Tianxing Chen, Jiaqi Liang, Zishun Shen, Haoran Lu,
423 Ruihai Wu, and Hao Dong. Garmentpile++: Affordance-driven cluttered garments retrieval
424 with vision-language reasoning. Preprint, 2026. CorpusID: 286238271.
- 425 [12] Ruihai Wu, Haoran Lu, Yiyan Wang, Yubo Wang, and Hao Dong. Unigarmentmanip: A
426 unified framework for category-level garment manipulation via dense visual correspondence.
427 *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
428 16340–16350, 2024. URL <https://api.semanticscholar.org/CorpusID:269757227>.
- 429 [13] Open X-Embodiment Collaboration, Abby O’Neill, Abdul Rehman, Abhinav Gupta, Abhi-
430 ram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim
431 Gupta, Ajay Mandekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan,
432 Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait
433 Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin,
434 Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-
435 Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles
436 Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang,
437 Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel

- 438 Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter B uchler, Dinesh
439 Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu,
440 Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gaurav S.
441 Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn,
442 Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben
443 Amor, Henrik I Christensen, Hiroki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang,
444 Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jae-
445 hyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg,
446 Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jian-
447 lan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik,
448 Jo o Silv rio, Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador,
449 Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman,
450 Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento
451 Kawaharazuka, Kevin Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty
452 Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng,
453 Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam
454 Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius
455 Memmel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian
456 Du, Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minh Heo, Mohan Kumar
457 Srirama, Mohit Sharma, Moo Jin Kim, Muhammad Zubair Irshad, Naoaki Kanazawa, Nicklas
458 Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur
459 Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi,
460 Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano,
461 Pierre Sermanet, Pieter Abbeel, Priya Sundaesan, Qiuyu Chen, Quan Vuong, Rafael Rafailov,
462 Ran Tian, Ria Doshi, Roberto Mart'in-Mart'in, Rohan Bajjal, Rosario Scalise, Rose Hendrix,
463 Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan
464 Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar
465 Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun Xu, Siddhant
466 Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal,
467 Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale,
468 Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada,
469 Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z.
470 Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent
471 Vanhoucke, Vitor Guizilini, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu
472 Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li,
473 Yansong Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng
474 Zhu, Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho,
475 Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang,
476 Yunfan Jiang, Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo
477 Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and Zipeng Lin. Open X-Embodiment: Robotic
478 learning datasets and RT-X models. <https://arxiv.org/abs/2310.08864>, 2023.
- 479 [14] Alexander Khazatsky, Karl Pertsch, Suraj Nair, et al. Droid: A large-scale in-the-wild robot
480 manipulation dataset. [arXiv:2403.12945](https://arxiv.org/abs/2403.12945), 2024.
- 481 [15] Modi Shi, Yuxiang Lu, Huijie Wang, and Shaoze Yang. Open-sourcing go-1: The bitter lessons
482 of building vla systems at scale. <https://opendrive1ab.com/OpenGO1/>, September 2025.
483 Blog post.
- 484 [16] Modi Shi, Yuxiang Lu, Huijie Wang, Chengen Xie, and Qingwen Bu. Introducing agibot world
485 colosseum: A large-scale manipulation platform for scalable and intelligent embodied systems.
486 <https://opendrive1ab.com/Agibot-World/>, March 2025. Blog post.
- 487 [17] Guo Ye, Zexi Zhang, Xu Zhao, Shang Wu, Haoran Lu, Shihan Lu, and Han Liu. Learning to
488 feel the future: Dreamtactvla for contact-rich manipulation. *ArXiv*, abs/2512.23864, 2025. URL
489 <https://api.semanticscholar.org/CorpusID:284350273>.
- 490 [18] Haoran Geng, Feishi Wang, Songlin Wei, Yuyang Li, Bangjun Wang, Boshi An, Charlie Tianyue
491 Cheng, Haozhe Lou, Peihao Li, Yen-Jen Wang, Yutong Liang, Dylan Goetting, Chaoyi Xu,
492 Haozhe Chen, Yuxi Qian, Yiran Geng, Jiageng Mao, Weikang Wan, Mingtong Zhang, Jian-
493 gran Lyu, Siheng Zhao, Jiazhao Zhang, Jialiang Zhang, Chengyang Zhao, Haoran Lu, Yufei

- 494 Ding, Ran Gong, Yuran Wang, Yuxuan Kuang, Ruihai Wu, Baoxiong Jia, Carlo Sferrazza,
495 Hao Dong, Siyuan Huang, Yue Wang, Jitendra Malik, and Pieter Abbeel. Roboverse: To-
496 wards a unified platform, dataset and benchmark for scalable and generalizable robot learn-
497 ing. *ArXiv*, abs/2504.18904, 2025. URL [https://api.semanticscholar.org/CorpusID:
498 277741767](https://api.semanticscholar.org/CorpusID:277741767).
- 499 [19] Haoran Lu, Yitong Li, Ruihai Wu, Chuanruo Ning, Yan Shen, and Hao Dong. Unigarment:
500 A unified simulation and benchmark for garment manipulation. Preprint, 2024. CorpusID:
501 275782214.
- 502 [20] Haoran Lu, Ruihai Wu, Yitong Li, Sijie Li, Ziyu Zhu, Chuanruo Ning, Yan Shen, Longzan Luo,
503 Yuanpei Chen, and Hao Dong. Garmentlab: A unified simulation and benchmark for garment
504 manipulation. *ArXiv*, abs/2411.01200, 2024. URL [https://api.semanticscholar.org/
505 CorpusID:273811186](https://api.semanticscholar.org/CorpusID:273811186).
- 506 [21] Lightwheel. Lightwheel kitchen: 3d kitchen asset collection for nvidia isaac sim. GitHub repos-
507 itory, 2025. URL https://github.com/LightwheelAI/Lightwheel_Kitchen. Open-
508 source 3D assets under CC BY-NC 4.0.
- 509 [22] Mayank Mittal, Pascal Roth, James Tigue, Antoine Richard, Octi Zhang, Peter Du, Antonio
510 Serrano-Munoz, Xinjie Yao, René Zurbrügg, Nikita Rudin, et al. Isaac lab: A gpu-accelerated
511 simulation framework for multi-modal robot learning. *arXiv preprint arXiv:2511.04831*, 2025.
- 512 [23] Lightwheel Team. Lw-benchhub: Lightwheel’s end-to-end embodied ai simulation platform.
513 GitHub repository, 2025. URL https://github.com/lightwheel-ai/lw_benchhub.
- 514 [24] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-
515 Martín, Chen Wang, Gabrael Levine, Wensi Ai, Benjamin Jose Martinez, Hang Yin, Michael
516 Lingelbach, Minjune Hwang, Ayano Hiranaka, Sujay Garlanka, Arman Aydin, Sharon Lee,
517 Jiankai Sun, Mona Anvari, Manasi Sharma, Dhruva Bansal, Samuel Hunter, Kyu-Young
518 Kim, Alan Lou, Caleb R. Matthews, Ivan Villa-Renteria, Jerry Huayang Tang, Claire Tang,
519 Fei Xia, Yunzhu Li, Silvio Savarese, Hyowon Gweon, C. Karen Liu, Jiajun Wu, Fei-Fei
520 Li, and Salesforce AI Research. Behavior-1k: A human-centered, embodied ai benchmark
521 with 1,000 everyday activities and realistic simulation. *ArXiv*, abs/2403.09227, 2024. URL
522 <https://api.semanticscholar.org/CorpusID:268520018>.
- 523 [25] Peijun Tang, Shangjin Xie, Binyan Sun, Baifu Huang, Kuncheng Luo, Haotian Yang, Weiqi
524 Jin, and Jianan Wang. Mind to hand: Purposeful robotic control via embodied reasoning, 2025.
525 URL <https://arxiv.org/abs/2512.08580>.
- 526 [26] Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yani Ding, Zhigang Wang,
527 Jiayuan Gu, Bin Zhao, Dong Wang, and Xuelong Li. Spatialvla: Exploring spatial repre-
528 sentations for visual-language-action model. *ArXiv*, abs/2501.15830, 2025. URL <https://api.semanticscholar.org/CorpusID:275921131>.
- 530 [27] Yida Niu, Xinhai Chang, Xin Liu, Ziyuan Jiao, and Yixin Zhu. Scalable trajectory generation
531 for whole-body mobile manipulation. In *CVPR*, 2026.
- 532 [28] Abhay Deshpande, Maya Guru, Rose Hendrix, Snehal Jauhri, Ainaz Eftekhari, Rohun Tripathi,
533 Max Argus, Jordi Salvador, Haoquan Fang, Matthew Wallingford, Wilbert Pumacay, Yejin Kim,
534 Quinn Pfeifer, Ying-Chun Lee, Piper Wolters, Omar Rayyan, Mingtong Zhang, Jiafei Duan,
535 Karen Farley, Winson Han, Eli Vanderbilt, Dieter Fox, Ali Farhadi, Georgia Chalvatzaki, Dhruv
536 Shah, and Ranjay Krishna. Molmob0t: Large-scale simulation enables zero-shot manipulation,
537 2026. URL <https://arxiv.org/abs/2603.16861>.
- 538 [29] Kaixuan Wang, Tianxing Chen, Jiawei Liu, Honghao Su, Shaolong Zhu, Minxuan Wang, Zixuan
539 Li, Yue Chen, Huan-ang Gao, Yusen Qin, Jiawei Wang, Qixuan Zhang, Lan Xu, Jingyi Yu, Yao
540 Mu, and Ping Luo. Manitwin: Scaling data-generation-ready digital object dataset to 100k.
541 Preprint, 2026. CorpusID: 286583937.
- 542 [30] Baijun Chen, Weijie Wan, Tianxing Chen, Xianda Guo, Congsheng Xu, Yuanyang Qi, Haojie
543 Zhang, Longyan Wu, Tianling Xu, Zixuan Li, et al. Univtac: A unified simulation platform
544 for visuo-tactile manipulation data generation, learning, and benchmarking. *arXiv preprint
545 arXiv:2602.10093*, 2026.

- 546 [31] Tianxing Chen, Zanxin Chen, Baijun Chen, Zijian Cai, Yibin Liu, Qiwei Liang, Zixuan Li,
547 XianLian Lin, Yiheng Ge, Zhenyu Gu, Weiliang Deng, Yubin Guo, Tian Nian, Xuanbing Xie,
548 Qiangyu Chen, Kailun Su, Tianling Xu, Guodong Liu, Mengkang Hu, Huan ang Gao, Kaixuan
549 Wang, Zhixuan Liang, Yusen Qin, Xiaokang Yang, Ping Luo, and Yao Mu. Robotwin 2.0: A
550 scalable data generator and benchmark with strong domain randomization for robust bimanual
551 robotic manipulation. *ArXiv*, abs/2506.18088, 2025. URL <https://api.semanticscholar.org/CorpusID:279999539>.
- 553 [32] Tianxing Chen, Yuran Wang, Mingleyang Li, Yanjiao Qin, Haochen Shi, Zixuan Li, Yifan
554 Hu, Yingsheng J Zhang, Kaixuan Wang, Yue Chen, Hongchen Wang, Renjing Xu, Ruihai Wu,
555 Yao Mu, Yaodong Yang, Hao Dong, and Ping Luo. Rmbench: Memory-dependent robotic
556 manipulation benchmark with insights into policy design. *ArXiv*, abs/2603.01229, 2026. URL
557 <https://api.semanticscholar.org/CorpusID:286222901>.
- 558 [33] Yang Tian, Yuyin Yang, Yiman Xie, Zetao Cai, Xu Shi, Ning Gao, Hangxu Liu, Xuekun
559 Jiang, Zherui Qiu, Feng Yuan, Yaping Li, Ping Wang, Junhao Cai, Jia Zeng, Hao Dong,
560 and Jiangmiao Pang. Interndata-a1: Pioneering high-fidelity synthetic data for pre-training
561 generalist policy. *ArXiv*, abs/2511.16651, 2025. URL <https://api.semanticscholar.org/CorpusID:283110021>.
- 563 [34] Yufei Wang, Zhou Xian, Feng Chen, Tsun-Hsuan Wang, Yian Wang, Katerina Fragkiadaki,
564 Zackory Erickson, David Held, and Chuang Gan. Robogen: Towards unleashing infinite data
565 for automated robot learning via generative simulation, 2023.
- 566 [35] Wenbo Wang, Fangyun Wei, QiXiu Li, Xi Chen, Yaobo Liang, Chang Xu, Jiaolong Yang,
567 and Baining Guo. Mobilemanibench: Simplifying model verification for mobile manipulation.
568 *arXiv preprint arXiv:2602.05233*, 2026.
- 569 [36] Balakumar Sundaralingam, Adithyavairavan Murali, and Stan Birchfield. curobov2: Dynamics-
570 aware motion generation with depth-fused distance fields for high-dof robots, 2026.
- 571 [37] Balakumar Sundaralingam, Siva Kumar Sastry Hari, Adam Fishman, Caelan Garrett, Karl Van
572 Wyk, Valts Blukis, Alexander Millane, Helen Oleynikova, Ankur Handa, Fabio Ramos, Nathan
573 Ratliff, and Dieter Fox. curobo: Parallelized collision-free minimum-jerk robot motion genera-
574 tion, 2023. URL <https://arxiv.org/abs/2310.17274>.
- 575 [38] Jiayi Chen, Yubin Ke, and He Wang. Bodex: Scalable and efficient robotic dexterous grasp
576 synthesis using bilevel optimization, 2025. URL <https://arxiv.org/abs/2412.16490>.
- 577 [39] Qingwei Ben, Feiyu Jia, Jia Zeng, Juntong Dong, Dahua Lin, and Jiangmiao Pang.
578 Homie: Humanoid loco-manipulation with isomorphic exoskeleton cockpit. *arXiv preprint*
579 *arXiv:2502.13013*, 2025.
- 580 [40] Zhengyi Luo, Ye Yuan, Tingwu Wang, Chenran Li, Sirui Chen, Fernando Castañeda, Zi-Ang
581 Cao, Jiefeng Li, David Minor, Qingwei Ben, Xingye Da, Runyu Ding, Cyrus Hogg, Lina Song,
582 Edy Lim, Eugene Jeong, Tairan He, Haoru Xue, Wenli Xiao, Zi Wang, Simon Yuen, Jan Kautz,
583 Yan Chang, Umar Iqbal, Linxi Fan, and Yuke Zhu. Sonic: Supersizing motion tracking for
584 natural humanoid whole-body control. *arXiv preprint arXiv:2511.07820*, 2025.
- 585 [41] Chengkai Hou, Zhengrong Xue, Bingyang Zhou, Jinghan Ke, Lin Shao, and Huazhe Xu.
586 Key-grid: Unsupervised 3d keypoints detection using grid heatmap features, 2024. URL
587 <https://arxiv.org/abs/2410.02237>.
- 588 [42] Ruoxi Shi, Zhengrong Xue, Yang You, and Cewu Lu. Skeleton merger: an unsupervised aligned
589 keypoint detector. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
590 *recognition*, pages 43–52, 2021.
- 591 [43] Chuanruo Ning, Ruihai Wu, Haoran Lu, Kaichun Mo, and Hao Dong. Where2explore: Few-shot
592 affordance learning for unseen novel categories of articulated objects. *ArXiv*, abs/2309.07473,
593 2023. URL <https://api.semanticscholar.org/CorpusID:261823492>.
- 594 [44] InternData-N1 Dataset contributors. Interndata-n1 dataset. <https://huggingface.co/datasets/InternRobotics/InternData-N1>, 2025.
- 595

- 596 [45] Abhay Deshpande, Maya Guru, Rose Hendrix, Snehal Jauhri, Haoquan Fang, Wilbert Pumacay,
597 Yejin Kim, Quinn Pfeifer, Ying-Chun Lee, Piper Wolters, Omar Rayyan, Mingtong Zhang,
598 Karen Farley, Winson Han, Eli Vanderbilt, Dieter Fox, Ali Farhadi, Georgia Chalvatzaki, Dhruv
599 Shah, and Ranjay Krishna. Molmob0t: Large-scale simulation enables zero-shot manipulation.
600 Preprint, 2026. CorpusID: 286498371.
- 601 [46] Ning Gao, Yilun Chen, Shuai Yang, Xinyi Chen, Yang Tian, Hao Li, Haifeng Huang, Hanqing
602 Wang, Tai Wang, and Jiangmiao Pang. Genmanip: Llm-driven simulation for generalizable
603 instruction-following manipulation. In *Proceedings of the Computer Vision and Pattern Recog-
604 nition Conference*, pages 12187–12198, 2025.
- 605 [47] Mingjie Pan, Jiyao Zhang, Tianshu Wu, Yinghao Zhao, Wenlong Gao, and Hao Dong. Omni-
606 manip: Towards general robotic manipulation via object-centric interaction primitives as spatial
607 constraints. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
608 Recognition (CVPR)*, pages 17359–17369, June 2025.
- 609 [48] Chenghao Yin, Da Huang, Di Yang, Jichao Wang, Nanshu Zhao, Chen Xu, Wenjun Sun, Linjie
610 Hou, Zhijun Li, Junhui Wu, Zhaobo Liu, Zhen Xiao, Sheng Zhang, Lei Bao, Rui Feng, Zhenquan
611 Pang, Jiayu Li, Qian Wang, and Maoqing Yao. Genie sim 3.0 : A high-fidelity comprehensive
612 simulation platform for humanoid robot, 2026. URL <https://arxiv.org/abs/2601.02078>.
- 613 [49] Yao Mu, Tianxing Chen, Zaxin Chen, Shijia Peng, Zhiqian Lan, Zeyu Gao, Zhixuan Liang,
614 Qiaojun Yu, Yude Zou, Mingkun Xu, Lunkai Lin, Zhiqiang Xie, Mingyu Ding, and Ping Luo.
615 Robotwin: Dual-arm robot benchmark with generative digital twins. In *Proceedings of the
616 Computer Vision and Pattern Recognition Conference (CVPR)*, pages 27649–27660, June 2025.
- 617 [50] Sizhe Yang, Yiman Xie, Zhixuan Liang, Yang Tian, Jia Zeng, Dahua Lin, and Jiangmiao Pang.
618 Ultralexgrasp: Learning universal dexterous grasping for bimanual robots with synthetic data,
619 2026. URL <https://arxiv.org/abs/2603.05312>.
- 620 [51] Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang.
621 Unidexgrasp++: Improving dexterous grasping policy learning via geometry-aware curriculum
622 and iterative generalist-specialist learning. In *Proceedings of the IEEE/CVF International
623 Conference on Computer Vision*, pages 3891–3902, 2023.
- 624 [52] Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzhen Xu, Puhao Li, Tengyu Liu, and He Wang.
625 Dexgraspnet: A large-scale robotic dexterous grasp dataset for general objects based on sim-
626 ulation. *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages
627 11359–11366, 2022. URL <https://api.semanticscholar.org/CorpusID:252734719>.
- 628 [53] Chen Bao, Helin Xu, Yuzhe Qin, and Xiaolong Wang. Dexart: Benchmarking generalizable
629 dexterous manipulation with articulated objects. In *Proceedings of the IEEE/CVF Conference
630 on Computer Vision and Pattern Recognition*, pages 21190–21200, 2023.
- 631 [54] Mu Lin, Yi-Lin Wei, Jiaxuan Chen, Yuhao Lin, Shuoyu Chen, Jiangran Lyu, Jiayi Chen, Yansong
632 Tang, He Wang, and Wei-Shi Zheng. Bidexgrasp: Coordinated bimanual dexterous grasps
633 across object geometries and sizes, 2026. URL <https://arxiv.org/abs/2604.06589>.
- 634 [55] Yufei Wang, Ziyu Wang, Mino Nakura, Pratik Bhowal, Chia-Liang Kuo, Yi-Ting Chen,
635 Zackory Erickson, and David Held. Articubot: Learning universal articulated object ma-
636 nipulation policy via large scale simulation. *ArXiv*, abs/2503.03045, 2025. URL <https://api.semanticscholar.org/CorpusID:276781877>.
- 637 [56] Qwen Team. Qwen3.5-omni technical report, 2026. URL <https://arxiv.org/abs/2604.15804>.
- 640 [57] Changfeng Ma, Yang Li, Xinhao Yan, Jiachen Xu, Yunhan Yang, Chunshi Wang, Zibo Zhao,
641 Yanwen Guo, Zhuo Chen, and Chunchao Guo. P3-sam: Native 3d part segmentation. *arXiv
642 preprint arXiv:2509.06784*, 2025.
- 643 [58] Xinhao Yan, Jiachen Xu, Yang Li, Changfeng Ma, Yunhan Yang, Chunshi Wang, Zibo Zhao,
644 Zeqiang Lai, Yunfei Zhao, Zhuo Chen, et al. X-part: High fidelity and structure coherent shape
645 decomposition. *arXiv preprint arXiv:2509.08643*, 2025.

- 646 [59] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. Rlbench: The robot
647 learning benchmark & learning environment. *arXiv preprint arXiv:1909.12271*, 2019. URL
648 <https://arxiv.org/abs/1909.12271>.
- 649 [60] Ajay Mandlekar, Soroush Nasiriany, Bowen Wen, Iretiayo Akinola, Yashraj Narang, Linxi Fan,
650 Yuke Zhu, and Dieter Fox. Mimicgen: A data generation system for scalable robot learning
651 using human demonstrations. In *Proceedings of The 7th Conference on Robot Learning*, volume
652 229 of *Proceedings of Machine Learning Research*, pages 1820–1864. PMLR, 2023. URL
653 <https://proceedings.mlr.press/v229/mandlekar23a.html>.
- 654 [61] Yufei Wang, Zhou Xian, Feng Chen, Tsun-Hsuan Wang, Yian Wang, Katerina Fragkiadaki,
655 Zackory Erickson, David Held, and Chuang Gan. Robogen: Towards unleashing infinite data
656 for automated robot learning via generative simulation. In *Proceedings of the 41st Interna-*
657 *tional Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning*
658 *Research*, pages 51936–51983. PMLR, 2024. URL <https://proceedings.mlr.press/v235/wang24cc.html>.
- 660 [62] Lirui Wang, Yiyang Ling, Zhecheng Yuan, Mohit Shridhar, Chen Bao, Yuzhe Qin, Bailin
661 Wang, Huazhe Xu, and Xiaolong Wang. Gensim: Generating robotic simulation tasks via large
662 language models. In *The Twelfth International Conference on Learning Representations*, 2024.
663 URL <https://openreview.net/forum?id=0I3RoHoWAN>.
- 664 [63] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Avnish Narayan, Hayden Shively,
665 Adithya Bellathur, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark
666 and evaluation for multi-task and meta reinforcement learning. *arXiv preprint arXiv:1910.10897*,
667 2019. URL <https://arxiv.org/abs/1910.10897>.
- 668 [64] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles
669 Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac
670 gym: High performance gpu-based physics simulation for robot learning. In *Thirty-Fifth*
671 *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.
672 URL <https://arxiv.org/abs/2108.10470>.
- 673 [65] OpenAI, Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Józefowicz, Bob McGrew,
674 Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, Jonas Schneider,
675 Szymon Sidor, Josh Tobin, Peter Welinder, Lilian Weng, and Wojciech Zaremba. Learning
676 dexterous in-hand manipulation. *arXiv preprint arXiv:1808.00177*, 2018. URL [https://](https://arxiv.org/abs/1808.00177)
677 arxiv.org/abs/1808.00177.
- 678 [66] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan
679 Aparicio Ojea, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with
680 synthetic point clouds and analytic grasp metrics. In *Robotics: Science and Systems*, 2017.
681 URL <https://www.roboticsproceedings.org/rss13/p58.pdf>.
- 682 [67] Clemens Eppner, Arsalan Mousavian, and Dieter Fox. Acronym: A large-scale grasp dataset
683 based on simulation. In *2021 IEEE International Conference on Robotics and Automa-*
684 *tion (ICRA)*. IEEE, 2021. URL [https://research.nvidia.com/publication/2021-05_](https://research.nvidia.com/publication/2021-05_acronym-large-scale-grasp-dataset-based-simulation)
685 [acronym-large-scale-grasp-dataset-based-simulation](https://research.nvidia.com/publication/2021-05_acronym-large-scale-grasp-dataset-based-simulation).
- 686 [68] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof graspnet: Variational grasp
687 generation for object manipulation. In *Proceedings of the IEEE/CVF International Conference*
688 *on Computer Vision*, pages 2901–2910, 2019. URL [https://openaccess.thecvf.](https://openaccess.thecvf.com/content_ICCV_2019/html/Mousavian_6-DOF_GraspNet_Variational_Grasp_Generation_for_Object_Manipulation_ICCV_2019_paper.html)
689 [com/content_ICCV_2019/html/Mousavian_6-DOF_GraspNet_Variational_Grasp_](https://openaccess.thecvf.com/content_ICCV_2019/html/Mousavian_6-DOF_GraspNet_Variational_Grasp_Generation_for_Object_Manipulation_ICCV_2019_paper.html)
690 [Generation_for_Object_Manipulation_ICCV_2019_paper.html](https://openaccess.thecvf.com/content_ICCV_2019/html/Mousavian_6-DOF_GraspNet_Variational_Grasp_Generation_for_Object_Manipulation_ICCV_2019_paper.html).
- 691 [69] Martin Sundermeyer, Arsalan Mousavian, Rudolph Triebel, and Dieter Fox. Contact-graspnet:
692 Efficient 6-dof grasp generation in cluttered scenes. *arXiv preprint arXiv:2103.14127*, 2021.
693 URL <https://arxiv.org/abs/2103.14127>.
- 694 [70] Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzhen Xu, Puhao Li, Tengyu Liu, and He Wang.
695 Dexgraspnet: A large-scale robotic dexterous grasp dataset for general objects based on sim-
696 ulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages
697 11359–11366. IEEE, 2023. URL <https://arxiv.org/abs/2210.02697>.

- 698 [71] Yanming Shao and Chenxi Xiao. Bimanual grasp synthesis for dexterous robot hands. *arXiv preprint arXiv:2411.15903*, 2024. URL <https://arxiv.org/abs/2411.15903>.
699
- 700 [72] Kaichun Mo, Leonidas Guibas, Mustafa Mukadam, Abhinav Gupta, and Shubham Tulsiani. Where2act: From pixels to actions for articulated 3d objects. *arXiv preprint arXiv:2101.02692*,
701 2021. URL <https://arxiv.org/abs/2101.02692>.
702
- 703 [73] Ruihai Wu, Yan Zhao, Kaichun Mo, Zizheng Guo, Yian Wang, Tianhao Wu, Qingnan Fan,
704 Xuelin Chen, Leonidas Guibas, and Hao Dong. Vat-mart: Learning visual action trajectory
705 proposals for manipulating 3d articulated objects. In *The Tenth International Conference on*
706 *Learning Representations*, 2022. URL <https://openreview.net/forum?id=iEx3PiooLy>.
- 707 [74] Haoran Geng, Helin Xu, Chengyang Zhao, Chao Xu, Li Yi, Siyuan Huang, and He Wang. Gapartnet: Cross-category domain-generalizable object perception and manipulation
708 via generalizable and actionable parts. In *Proceedings of the IEEE/CVF Confer-*
709 *ence on Computer Vision and Pattern Recognition*, pages 7081–7091, 2023. URL
710 [https://openaccess.thecvf.com/content/CVPR2023/html/Geng_GAPartNet_](https://openaccess.thecvf.com/content/CVPR2023/html/Geng_GAPartNet_Cross-Category_Domain-Generalizable_Object_Perception_and_Manipulation_via_Generalizable_and_CVPR_2023_paper.html)
711 [Cross-Category_Domain-Generalizable_Object_Perception_and_Manipulation_](https://openaccess.thecvf.com/content/CVPR2023/html/Geng_GAPartNet_Cross-Category_Domain-Generalizable_Object_Perception_and_Manipulation_via_Generalizable_and_CVPR_2023_paper.html)
712 [via_Generalizable_and_CVPR_2023_paper.html](https://openaccess.thecvf.com/content/CVPR2023/html/Geng_GAPartNet_Cross-Category_Domain-Generalizable_Object_Perception_and_Manipulation_via_Generalizable_and_CVPR_2023_paper.html).
713
- 714 [75] Kaichun Mo, Yuzhe Qin, Fanbo Xiang, Hao Su, and Leonidas Guibas. O2o-afford: Annotation-
715 free large-scale object-object affordance learning. In *Conference on Robot Learning*, 2021.
716 URL <https://openreview.net/forum?id=EougVeukEH9>.
- 717 [76] Yifan You, Lin Shao, Toki Migimatsu, and Jeannette Bohg. Omnihang: Learning to hang
718 arbitrary objects using contact point correspondences and neural collision estimation. In
719 *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021. URL
720 <https://arxiv.org/abs/2103.14283>.
- 721 [77] Bingjie Tang, Michael A. Lin, Ireteayo Akinola, Ankur Handa, Gaurav S. Sukhatme, Fabio
722 Ramos, Dieter Fox, and Yashraj Narang. Industreal: Transferring contact-rich assembly
723 tasks from simulation to reality. In *Robotics: Science and Systems*, 2023. URL <https://arxiv.org/abs/2305.17110>.
724
- 725 [78] Shengheng Deng, Xun Xu, Chaozheng Wu, Ke Chen, and Kui Jia. 3d affordancenet:
726 A benchmark for visual object affordance understanding. In *Proceedings of the*
727 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1778–1787,
728 2021. URL [https://openaccess.thecvf.com/content/CVPR2021/html/Deng_3D_](https://openaccess.thecvf.com/content/CVPR2021/html/Deng_3D_AffordanceNet_A_Benchmark_for_Visual_Object_Affordance_Understanding_CVPR_2021_paper.html)
729 [AffordanceNet_A_Benchmark_for_Visual_Object_Affordance_Understanding_](https://openaccess.thecvf.com/content/CVPR2021/html/Deng_3D_AffordanceNet_A_Benchmark_for_Visual_Object_Affordance_Understanding_CVPR_2021_paper.html)
730 [CVPR_2021_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Deng_3D_AffordanceNet_A_Benchmark_for_Visual_Object_Affordance_Understanding_CVPR_2021_paper.html).
- 731 [79] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo
732 Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
733 URL <https://arxiv.org/abs/1512.03012>.
734
- 735 [80] Kaichun Mo, Shilin Zhu, Angel X. Chang, Li Yi, Subarna Tripathi, Leonidas J. Guibas, and
736 Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object
737 understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
738 *Recognition*, pages 909–918, 2019. URL [https://openaccess.thecvf.com/content_](https://openaccess.thecvf.com/content_CVPR_2019/html/Mo_PartNet_A_Large-Scale_Benchmark_for_Fine-Grained_and_Hierarchical_Part-Level_3D_CVPR_2019_paper.html)
739 [CVPR_2019/html/Mo_PartNet_A_Large-Scale_Benchmark_for_Fine-Grained_and_](https://openaccess.thecvf.com/content_CVPR_2019/html/Mo_PartNet_A_Large-Scale_Benchmark_for_Fine-Grained_and_Hierarchical_Part-Level_3D_CVPR_2019_paper.html)
740 [Hierarchical_Part-Level_3D_CVPR_2019_paper.html](https://openaccess.thecvf.com/content_CVPR_2019/html/Mo_PartNet_A_Large-Scale_Benchmark_for_Fine-Grained_and_Hierarchical_Part-Level_3D_CVPR_2019_paper.html).
- 741 [81] SAPIEN. Partnet-mobility dataset. <https://sapien.ucsd.edu/>, 2020. Official dataset entry
742 point within the SAPIEN ecosystem.
- 743 [82] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua
744 Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and
745 Hao Su. Sapien: A simulated part-based interactive environment. In *Proceedings of the*
746 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11097–11107, 2020.
747 URL [https://openaccess.thecvf.com/content_CVPR_2020/html/Xiang_SAPIEN_](https://openaccess.thecvf.com/content_CVPR_2020/html/Xiang_SAPIEN_A_Simulated_Part-Based_Interactive_Environment_CVPR_2020_paper.html)
748 [A_Simulated_Part-Based_Interactive_Environment_CVPR_2020_paper.html](https://openaccess.thecvf.com/content_CVPR_2020/html/Xiang_SAPIEN_A_Simulated_Part-Based_Interactive_Environment_CVPR_2020_paper.html).

- 749 [83] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli Vander-
750 Bilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Obja-
751 verse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Confer-*
752 *ence on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. URL
753 [https://openaccess.thecvf.com/content/CVPR2023/html/Deitke_Objaverse_A_](https://openaccess.thecvf.com/content/CVPR2023/html/Deitke_Objaverse_A_Universe_of_Annotated_3D_Objects_CVPR_2023_paper.html)
754 [Universe_of_Annotated_3D_Objects_CVPR_2023_paper.html](https://openaccess.thecvf.com/content/CVPR2023/html/Deitke_Objaverse_A_Universe_of_Annotated_3D_Objects_CVPR_2023_paper.html).
- 755 [84] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using
756 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. URL [https://arxiv.org/abs/2209.](https://arxiv.org/abs/2209.14988)
757 [14988](https://arxiv.org/abs/2209.14988).
- 758 [85] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv*
759 *preprint arXiv:2305.02463*, 2023. URL <https://arxiv.org/abs/2305.02463>.
- 760 [86] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based
761 control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages
762 5026–5033. IEEE, 2012. doi: 10.1109/IROS.2012.6386109.
- 763 [87] Benjamin Ellenberger. Pybullet gymperium. [https://github.com/benelot/](https://github.com/benelot/pybullet-gym)
764 [pybullet-gym](https://github.com/benelot/pybullet-gym), 2018–2019.
- 765 [88] NVIDIA. *Isaac Sim Documentation*, 2025. URL [https://docs.isaacsim.omniverse.](https://docs.isaacsim.omniverse.nvidia.com/4.5.0/index.html)
766 [nvidia.com/4.5.0/index.html](https://docs.isaacsim.omniverse.nvidia.com/4.5.0/index.html). Official documentation.
- 767 [89] Yuke Zhu, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Abhishek Joshi, Kevin Lin,
768 Abhiram Maddukuri, Soroush Nasiriany, and Yifeng Zhu. robosuite: A modular simulation
769 framework and benchmark for robot learning. *arXiv preprint arXiv:2009.12293*, 2020. URL
770 <https://arxiv.org/abs/2009.12293>.
- 771 [90] Tongzhou Mu, Zhan Ling, Fanbo Xiang, Derek Yang, Xuanlin Li, Stone Tao, Zhiao Huang,
772 Zhiwei Jia, and Hao Su. Maniskill: Generalizable manipulation skill benchmark with large-scale
773 demonstrations. *arXiv preprint arXiv:2107.14483*, 2021. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2107.14483)
774 [2107.14483](https://arxiv.org/abs/2107.14483).
- 775 [91] Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang, Stone
776 Tao, Xinyue Wei, Yunchao Yao, Xiaodi Yuan, Pengwei Xie, Zhiao Huang, Rui Chen, and Hao
777 Su. Maniskill2: A unified benchmark for generalizable manipulation skills. In *International*
778 *Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=b_CQDy9vrD1)
779 [id=b_CQDy9vrD1](https://openreview.net/forum?id=b_CQDy9vrD1).
- 780 [92] Sanjana Srivastava, Chengshu Li, Michael Lingelbach, Roberto Martín-Martín, Fei Xia, Kent
781 Vainio, Zheng Lian, Cem Gokmen, Shyamal Buch, C. Karen Liu, Silvio Savarese, Hyowon
782 Gweon, Jiajun Wu, and Li Fei-Fei. Behavior: Benchmark for everyday household activities in
783 virtual, interactive, and ecological environments. *arXiv preprint arXiv:2108.03332*, 2021. URL
784 <https://arxiv.org/abs/2108.03332>.
- 785 [93] Richard S. Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A
786 framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1-2):
787 181–211, 1999. doi: 10.1016/S0004-3702(99)00052-1.
- 788 [94] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David,
789 Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine
790 Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey,
791 Sally Jesmonth, Nikhil J. Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei
792 Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao,
793 Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan,
794 Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan,
795 and Andy Zeng. Do as i can, not as i say: Grounding language in robotic affordances. In *Confer-*
796 *ence on Robot Learning*, 2022. URL https://openreview.net/forum?id=bdHkMjBJG_w.
- 797 [95] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence,
798 and Andy Zeng. Code as policies: Language model programs for embodied control. *arXiv*
799 *preprint arXiv:2209.07753*, 2022. URL <https://arxiv.org/abs/2209.07753>.

800 [96] Weiyu Liu, Jiayuan Mao, Joy Hsu, Tucker Hermans, Animesh Garg, and Jiajun Wu. Composable
801 part-based manipulation. In *Proceedings of The 7th Conference on Robot Learning*, volume
802 229 of *Proceedings of Machine Learning Research*, pages 1300–1315. PMLR, 2023. URL
803 <https://proceedings.mlr.press/v229/liu23e.html>.

804 A Language Annotation Pipeline

805 A.1 Asset-level Language Annotation Details

806 **Automatic multi-view rendering.** For each object-level asset, we automatically render multi-view
807 observations before querying the VLM. We first compute the asset’s axis-aligned bounding box and
808 use its center and diagonal length to normalize camera placement. Specifically, let c denote the
809 bounding-box center and d denote the diagonal length. We place eight cameras on a horizontal circle
810 around the object with evenly spaced azimuth angles $\{0^\circ, 45^\circ, 90^\circ, 135^\circ, 180^\circ, 225^\circ, 270^\circ, 315^\circ\}$,
811 and set the camera distance proportional to d so that the object is centered and consistently scaled in
812 every rendered image. All cameras look at c and render RGB observations in simulation. We use eight
813 views by default because they provide sufficient coverage of most object geometries while remaining
814 close to the practical multi-image input budget for stable VLM reasoning in our implementation.

815 **Global hierarchical description.** Given the eight rendered views, we prompt Qwen3.5-VL [56]
816 to generate asset-level language annotations at three granularities: a phrase-level label, a sentence-
817 level description, and a paragraph-level description. The phrase-level label provides a compact
818 semantic index of the object, the sentence-level description summarizes its main function and major
819 manipulable components, and the paragraph-level description provides dense interaction reasoning
820 about object parts, affordances, feasible actions, and manipulation constraints.

821 **Part-aware dense description.** To make the dense paragraph grounded in object structure, we
822 further query the VLM with part-specific masked renderings. For each segmented part, we render
823 an image where the target part is highlighted while the remaining object is kept as context. The
824 VLM is then asked to describe only the highlighted part, including its part name, visual evidence,
825 function, feasible robot interactions, and manipulation constraints. Each part-level description is
826 associated with its segmentation ID, which allows the final dense language annotation to be traced
827 back to concrete visual regions rather than free-form text alone.

828 **Output format and filtering.** We require the VLM to produce a structured JSON output. To
829 reduce hallucination, the prompt explicitly asks the model to describe only visible evidence, avoid
830 unsupported assumptions, and mark uncertain attributes as unknown. We discard outputs that cannot
831 be parsed or whose part descriptions cannot be matched to a valid segmentation ID. The resulting
832 language annotations are used as semantic priors for visual grounding and physics-based action
833 generation.

834 **Prompt for global asset-level language annotation.** We use the following prompt to generate
835 phrase-, sentence-, and paragraph-level asset descriptions from the eight rendered views.

Global asset-level language prompt

```
System: You are a careful annotation engine for robot manipulation assets.  
Your task is to describe the visible 3D object for robot reasoning and  
manipulation. Use only visual evidence from the provided views. Do not invent  
hidden mechanisms or unsupported functions. If an attribute is unclear, write  
"unknown".  
User: You are given eight RGB renderings of the same 3D object from surrounding  
viewpoints. Please analyze the object as a robot manipulation asset and output  
a structured annotation with three levels of detail.  
Level 1: phrase-level label. Write a short phrase, no more than 6 words, that  
identifies the object category or functional identity.  
Level 2: sentence-level description. Write one sentence, no more than 30  
words, describing the object’s main function and major manipulable parts.
```

836

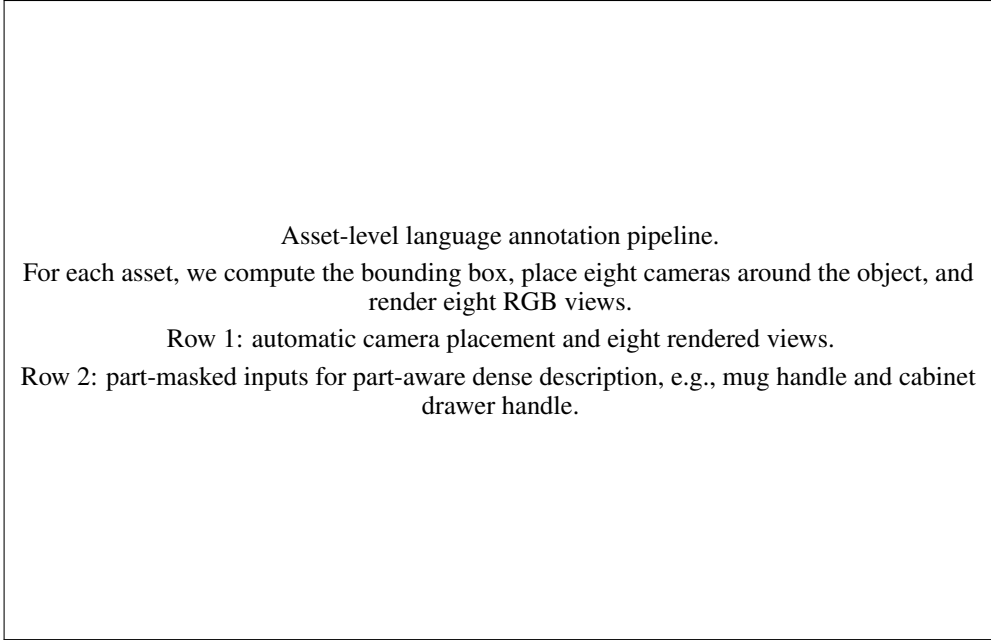


Figure 6: Supplementary illustration of asset-level language annotation. We automatically normalize camera placement using the object bounding box and render eight surrounding views for global asset-level description. For paragraph-level dense reasoning, segmented parts are individually highlighted and fed to the VLM, enabling part-aware descriptions grounded in visual evidence.

Level 3: paragraph-level description. Write one dense paragraph for robot manipulation reasoning. Describe the object’s visible parts, their functions, likely affordances, feasible interactions, and manipulation constraints. Focus on what a robot can grasp, pull, push, support, open, insert into, hang on, or avoid.

Return the result in JSON format with the following fields:
 phrase, sentence, paragraph, major_parts, uncertain_attributes.

837

838 **Prompt for part-aware dense description.** For paragraph-level annotation, we further query the
 839 VLM with masked part renderings.

Part-aware language prompt

System: You are a robot manipulation annotation engine. You will be given an image of a 3D object where one segmented part is highlighted. The highlighted region is the target part. The rest of the object is shown only as context. Describe only the highlighted part.

User: Please analyze the highlighted part for robot manipulation. Do not describe unrelated parts unless they are needed to explain the target part. If the highlighted part’s function is uncertain, write "unknown" instead of guessing.

Please answer the following questions: 1. What is the name of the highlighted part? 2. What visual evidence supports this part name? 3. What is the function of this part? 4. What robot interactions are feasible with this part? 5. What constraints should a robot consider when interacting with this part? 6. How is this part related to the whole object?

Return the result in JSON format with the following fields:
 part_name, visual_evidence, part_function, feasible_robot_interactions, manipulation_constraints, relation_to_whole_object, confidence.

840

841 **Example output: mug.**

Global asset-level output: mug

Phrase-level label: Blue handled mug.

Sentence-level description: A blue mug with a hollow cup body and a side handle, designed for containing liquid and being grasped by the handle.

Paragraph-level description: The object is a mug composed of a cylindrical cup body, an open rim, a hollow interior, a side handle, and a bottom base. The cup body provides the main container volume and can be used for placement and containment reasoning. The side handle is the most important manipulation part because it supports stable grasping, lifting, and carrying. The open rim and hollow interior indicate that the mug can contain liquid or small objects, but they should be avoided as forceful grasp targets. The bottom base supports stable placement on flat surfaces.

842

Part-aware output: mug handle

Highlighted part: Mug handle.

The highlighted part is the mug handle. It is a curved loop attached to the side of the cup body. Its function is to provide a safe and stable grasping region for lifting or carrying the mug. Feasible robot interactions include grasping the outer handle, inserting fingers or gripper tips through the loop, and pulling or lifting the mug through the handle. The robot should avoid colliding with the cup body near the handle attachment points and should align the grasp with the handle opening.

843

844 Example output: cabinet.

Global asset-level output: cabinet

Phrase-level label: Three-drawer cabinet.

Sentence-level description: A wooden cabinet with stacked front drawers, horizontal handles, side panels, and an open top cavity for storage and drawer manipulation.

Paragraph-level description: The object is a cabinet with multiple stacked drawer fronts, horizontal metal handles, side panels, an open top cavity, and a supporting base. The drawer fronts are likely movable storage components that translate outward when pulled. The handles are the primary interaction parts for grasping and pulling the drawers. The side panels and frame provide rigid structural support and should generally be treated as non-manipulable surfaces. The open top cavity can be used for placing or retrieving objects, while the base supports the cabinet on the floor.

845

Part-aware output: cabinet drawer handle

Highlighted part: Drawer handle.

The highlighted part is a drawer handle on the front of the cabinet. It is an elongated horizontal bar attached to a drawer panel and is designed as the main interaction target for opening the drawer. Feasible robot interactions include parallel-jaw grasping, dexterous hooking, and pulling along the drawer's translation direction. The robot should approach the handle from the front, avoid collisions with neighboring drawer panels, and keep the pulling motion aligned with the drawer direction.

846

847 A.2 Room-level Language Annotation Details

848 **Automatic room-view rendering.** For each room-scale scene, we generate room-level language
849 annotations from both egocentric room views and a top-view spatial representation. We first de-
850 compose the scene into individual rooms and randomly sample camera poses inside each room.
851 Since uniformly sampled cameras may face walls or empty regions, we use simulator ground-truth
852 segmentation as an automatic view-quality signal. For each candidate camera pose, we render an
853 RGB image and its segmentation map, count the number of visible non-structural object instances,
854 and reject views with too few visible objects or views dominated by walls and floors. We keep a fixed
855 number of valid views, eight in our implementation, with diverse positions and viewing directions.
856 This produces informative room-level observations without manual camera placement.

857 **Top-view layout input.** In addition to the selected RGB views, we provide the VLM with a top-
858 view occupancy map of the room or scene. The map marks free space, obstacles, room boundaries,

859 and object instances with their corresponding labels. The RGB views provide local appearance and
860 object-level evidence, while the labeled occupancy map provides global layout, spatial relations, and
861 navigable structure. This combination allows the VLM to reason about both what objects are present
862 and how they are arranged in the room.

863 **Hierarchical room-level description.** Given the selected room views and the labeled top-view
864 occupancy map, we prompt Qwen3.5-VL [56] to generate room-level language annotations at three
865 granularities. The phrase-level annotation summarizes the room type or global layout pattern. The
866 sentence-level annotation describes the major furniture, functional zones, and spatial arrangement.
867 The paragraph-level annotation provides dense scene-level reasoning, including object relations, navigable
868 regions, interaction-relevant areas, and possible long-horizon robot tasks. These descriptions
869 complement asset-level annotations by explaining how individual objects compose into a functional
870 scene.

871 **Output format and filtering.** We require structured JSON outputs and explicitly instruct the VLM
872 to use the labeled occupancy map for global layout reasoning and the RGB views for local visual
873 evidence. To reduce hallucination, the prompt asks the model to avoid unsupported object names or
874 hidden room functions and to mark uncertain attributes as unknown. We discard outputs that cannot
875 be parsed or whose object references cannot be matched to the labeled occupancy map.

876 **Prompt for room-level language annotation.** We use the following prompt to generate short-
877 sentence, long-sentence, and paragraph-level room descriptions from selected room views and the
878 labeled top-view occupancy map.

Room-level language prompt

System: You are a careful annotation engine for robot navigation and manipulation scenes. Your task is to describe the visible room or indoor scene for robot reasoning. Use the provided RGB views for local visual evidence and the labeled top-view occupancy map for global layout, spatial relations, and navigable structure. Do not invent hidden rooms, unsupported objects, or functions that are not visible or labeled. If an attribute is unclear, write "unknown".

User: You are given multiple RGB renderings sampled inside a room-scale scene, together with a top-view occupancy map. The occupancy map marks free space, obstacles, room boundaries, and object instances with labels. Please analyze the scene as a robot interaction environment and output a structured room-level annotation with three levels of detail.

Level 1: short-sentence description. Write one short sentence, no more than 15 words, summarizing the room type and global layout.

Level 2: long-sentence description. Write one to two long sentences describing the major furniture, functional zones, spatial arrangement, and navigation structure.

Level 3: paragraph-level scene description. Write one dense paragraph for robot navigation and manipulation reasoning. Describe the room layout, object relations, navigable regions, interaction-relevant areas, and possible long-horizon robot tasks. Focus on what a robot can navigate to, avoid, approach, pick, place, open, retrieve, or interact with.

Return the result in JSON format with the following fields: short_sentence, long_sentence, paragraph, room_type, functional_zones, object_inventory, spatial_relations, navigable_regions, interaction_relevant_areas, possible_robot_tasks, uncertain_attributes.

879

880 **Example output: multi-room scene.**

Room-level output: multi-room furnished apartment

Short-sentence description: A multi-room furnished apartment with connected rooms and open navigation corridors.

881

Long-sentence description: The scene contains several enclosed rooms arranged around central open spaces and corridors, with furniture mostly placed near walls. Major interaction areas include tables, cabinets, shelves, beds or sofas, and small objects on support surfaces.

Paragraph-level description: The scene contains multiple connected rooms with corridors, open passages, and several functional zones. The top-view layout shows central open areas surrounded by smaller rooms, with furniture arranged mostly along walls and corners. Visible interaction-relevant objects include tables, cabinets, shelves, beds or sofas, lamps, and small tabletop items. The main navigable regions are the open floor areas, corridors, and doorway passages, while dense furniture clusters and narrow passages should be treated as collision-sensitive regions. A robot can perform long-horizon tasks such as navigating between rooms, approaching furniture, retrieving objects from tables or shelves, placing objects on support surfaces, and interacting with storage furniture.

882

883 **Example output: dining room.**

Room-level output: dining room

Short-sentence description: A rectangular dining room with a central table and surrounding chairs.

Long-sentence description: The room contains a central dining area formed by a round table and multiple chairs, with additional support and storage furniture near the rear and side walls. Open floor space around the furniture provides navigation routes, while the table, chairs, and shelf create local obstacles and approach constraints.

Paragraph-level description: The room is a rectangular dining or meeting space with most interaction objects concentrated near the center and right wall. A round table with several chairs forms the main functional zone for placing, picking, and tabletop manipulation. A secondary table near the rear wall provides another support surface, while the side bookshelf can be used for object retrieval or placement. The open floor around the furniture provides navigable space, but the chairs and table legs create local obstacles that require careful path planning. A robot can navigate around the table, approach chairs or tables, retrieve objects from the table or shelf, and place objects on available support surfaces.

884

885 **B Visual Annotation Pipeline**

886 **B.1 Asset-level Visual Annotation Details**

887 **Multi-view point-cloud reconstruction.** For each object-level asset, we construct the input point
888 cloud from multi-view RGB-D observations rather than directly sampling points from the mesh.
889 This design makes the annotation pipeline closer to real-world inference, where object geometry is
890 reconstructed from camera observations. Similar to asset-level language annotation, we first compute
891 the object bounding box and place eight cameras around the asset with evenly spaced azimuth angles.
892 Each camera is oriented toward the bounding-box center and its distance is normalized by the object
893 scale.

894 Given the depth image D_i from camera i , camera intrinsics K_i , and camera-to-world extrinsics
895 $T_{c_i}^w = [R_i \mid t_i]$, each valid depth pixel (u, v) is back-projected into 3D as

$$\mathbf{x}_{uv}^{c_i} = D_i(u, v)K_i^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}, \quad \mathbf{x}_{uv}^w = R_i \mathbf{x}_{uv}^{c_i} + t_i.$$

896 The fused object point cloud is obtained by aggregating valid points from all rendered views:

$$\mathcal{P}_{cam} = \bigcup_{i=1}^N \{\mathbf{x}_{uv}^w \mid D_i(u, v) > 0, M_i(u, v) = 1\},$$

897 where $N = 8$ by default and M_i denotes the object foreground mask from simulation. We remove
898 background points using the mask and merge all points into the object coordinate frame.

899 **Coverage verification with mesh-based reference points.** Although the final annotation point
900 cloud is reconstructed from rendered observations, we use mesh-based sampling only as a coverage

901 check. Specifically, we sample a reference point set \mathcal{P}_{mesh} from the object mesh and compare it with
 902 the camera-fused point cloud \mathcal{P}_{cam} inside the object’s 3D region of interest. We measure the surface
 903 coverage ratio as

$$\text{Cov}(\mathcal{P}_{cam}, \mathcal{P}_{mesh}) = \frac{1}{|\mathcal{P}_{mesh}|} \sum_{\mathbf{p} \in \mathcal{P}_{mesh}} \mathbb{1} \left[\min_{\mathbf{q} \in \mathcal{P}_{cam}} \|\mathbf{p} - \mathbf{q}\|_2 < \epsilon \right],$$

904 where ϵ is a distance threshold. We also compare the reconstructed 3D region of interest with the
 905 mesh-level region using voxel occupancy:

$$\text{IoU}_{roi} = \frac{|\mathcal{V}(\mathcal{P}_{cam}) \cap \mathcal{V}(\mathcal{P}_{mesh})|}{|\mathcal{V}(\mathcal{P}_{cam}) \cup \mathcal{V}(\mathcal{P}_{mesh})|}.$$

906 If either the surface coverage or ROI overlap is below a predefined threshold, we add extra cameras
 907 targeting poorly covered regions and repeat the reconstruction. This adaptive step ensures that thin
 908 structures, concave regions, and partially occluded parts are sufficiently observed while preserving
 909 the camera-based nature of the final point cloud.

910 **Point-cloud downsampling.** The fused point cloud can contain a large number of points, especially
 911 for high-resolution renderings. Before part segmentation, we downsample \mathcal{P}_{cam} to $M = 80,000$
 912 points using farthest point sampling (FPS). Starting from an initial point set \mathcal{S}_1 , FPS iteratively selects

$$\mathbf{p}^* = \arg \max_{\mathbf{p} \in \mathcal{P}_{cam}} \min_{\mathbf{s} \in \mathcal{S}} \|\mathbf{p} - \mathbf{s}\|_2,$$

913 and updates $\mathcal{S} \leftarrow \mathcal{S} \cup \{\mathbf{p}^*\}$ until $|\mathcal{S}| = M$. This preserves global shape coverage and avoids
 914 over-concentrating points on large planar surfaces. The downsampled point cloud is denoted as \mathcal{P}_{80k} .

915 **Native 3D part segmentation.** We perform asset-level part segmentation on \mathcal{P}_{80k} using the Hun-
 916 yuan3D part segmentation model with 182M parameters, built on native 3D part decomposition
 917 modules such as P3-SAM [57] and X-Part [58]. The model predicts a set of part masks

$$\mathcal{M} = \{m_k\}_{k=1}^K, \quad m_k \in \{0, 1\}^{|\mathcal{P}_{80k}|},$$

918 where each m_k corresponds to one segmented part. For downstream modules that require denser
 919 geometry, we propagate part labels from \mathcal{P}_{80k} back to the original fused cloud by nearest-neighbor
 920 assignment:

$$\ell(\mathbf{p}) = \ell \left(\arg \min_{\mathbf{q} \in \mathcal{P}_{80k}} \|\mathbf{p} - \mathbf{q}\|_2 \right), \quad \mathbf{p} \in \mathcal{P}_{cam}.$$

921 In our implementation, the segmentation stage takes approximately half a minute per object, making
 922 it practical for large-scale automatic annotation.

923 **Failure cases.** The main failure mode of asset-level part segmentation is semantic granularity
 924 mismatch. Some objects are segmented at a coarser level than required for manipulation. For
 925 example, a wine glass may be predicted as a single part, while functional manipulation would
 926 benefit from separating the bowl, stem, and base. Such over-merged parts can limit part-aware
 927 language annotation and functional grasp generation, because different subparts correspond to
 928 different interaction strategies. We therefore treat part segmentation as a visual prior rather than a
 929 final executable label; downstream physics-based validation and task-specific geometric reasoning
 930 are still required to produce executable action annotations.

931 B.2 Room-level Visual Annotation Details

932 **Wall-based room decomposition.** For each room-scale scene, we first decompose the scene into
 933 individual rooms using wall boundaries and room-layout structure. Each room is represented as a
 934 2D floor-plan polygon \mathcal{R}_j in the global scene coordinate frame. This decomposition allows us to
 935 scan and annotate each room independently, while later stitching the resulting maps into a global
 936 room-level representation.

937 **Rule-based scan-pose sampling.** For each room \mathcal{R}_j , we randomly sample candidate camera and
 938 LiDAR scan poses inside the room polygon. A sampled pose is accepted only if it satisfies simple
 939 rule-based constraints: it must lie in free space, maintain a minimum distance from walls and large
 940 obstacles, and have a valid viewing direction toward the room interior rather than directly facing a
 941 nearby wall. We use simulator ground-truth segmentation only as an automatic view-quality signal.
 942 For each candidate pose π , we render an RGB image and an instance segmentation map, and compute
 943 the set of visible object IDs:

$$\mathcal{V}(\pi) = \left\{ o \in \mathcal{O}_j \mid \sum_{u,v} \mathbf{1}[S_\pi(u,v) = o] > \tau_{\text{area}} \right\},$$

944 where \mathcal{O}_j is the set of object instances in room \mathcal{R}_j , S_π is the rendered segmentation map, and τ_{area}
 945 filters out tiny or barely visible objects. We then greedily select scan poses that maximize object-ID
 946 coverage:

$$\pi_t = \arg \max_{\pi \in \Pi_j} \left| \mathcal{V}(\pi) \setminus \bigcup_{k < t} \mathcal{V}(\pi_k) \right|,$$

947 until the coverage ratio

$$\text{Cov}_{obj} = \frac{|\bigcup_t \mathcal{V}(\pi_t)|}{|\mathcal{O}_j|}$$

948 exceeds a predefined threshold or a maximum number of poses is reached. This prevents the room-
 949 level annotation pipeline from using uninformative views that mostly face walls or empty regions,
 950 and ensures that nearly all object instances in each room are observed at least once.

951 **Multi-height occupancy from LiDAR and ray tracing.** Given the selected scan poses, we
 952 construct room-level occupancy maps using simulated LiDAR and ray tracing. For each pose π_t and
 953 each height level $h_\ell \in \mathcal{H}$, we cast horizontal rays in uniformly sampled directions. For a ray direction
 954 \mathbf{d}_α , let ρ_α be the first-hit distance returned by ray tracing. The cells before the hit are marked as free,
 955 while the hit cell is marked as occupied:

$$\mathcal{F}_{t,\ell,\alpha} = \{g(\mathbf{x}_t + r\mathbf{d}_\alpha, h_\ell) \mid 0 < r < \rho_\alpha - \delta\}, \quad \mathcal{O}_{t,\ell,\alpha} = g(\mathbf{x}_t + \rho_\alpha\mathbf{d}_\alpha, h_\ell),$$

956 where \mathbf{x}_t is the scan position, $g(\cdot)$ maps a 3D point to a 2D grid cell, and δ is a small safety margin
 957 before the surface hit. We aggregate all rays with a log-odds occupancy update:

$$L_\ell(c) \leftarrow L_\ell(c) + \begin{cases} \lambda_{\text{occ}}, & c \in \mathcal{O}_{t,\ell,\alpha}, \\ \lambda_{\text{free}}, & c \in \mathcal{F}_{t,\ell,\alpha}, \\ 0, & \text{otherwise,} \end{cases}$$

958 and obtain the occupancy probability by $P_\ell(c) = \sigma(L_\ell(c))$. The final multi-height occupancy
 959 annotation is

$$\mathcal{G}_j = \{G_j^{h_1}, G_j^{h_2}, \dots, G_j^{h_L}\},$$

960 where different height slices capture floor-level free space, low obstacles, furniture geometry, and
 961 wall-level structure.

962 **Global map stitching.** After scanning each room independently, we stitch all room-level maps into
 963 a global scene representation using the floor-plan alignment and the global coordinate frame. Let T_j^w
 964 denote the transform from the local room frame to the global world frame. For each room map G_j ,
 965 we project its grid cells into the global map:

$$G^w(c) = \text{Fuse}_j(G_j((T_j^w)^{-1}c)),$$

966 where $\text{Fuse}(\cdot)$ combines overlapping evidence from adjacent rooms and shared boundaries. This
 967 produces a globally consistent occupancy map, object-layout map, and top-view representation for
 968 the entire scene.

969 **VLM-assisted wall and floor-plan filtering.** Raw occupancy maps often contain both structural
 970 elements and movable objects such as tables, chairs, cabinets, and shelves. To obtain a clean floor
 971 plan and wall-structure annotation, we use the VLM as a structural component selector. We first
 972 generate a top-view map with connected components and object labels, then prompt the VLM to keep
 973 persistent room-structure components, such as walls, doors, and room boundaries, while filtering out
 974 movable furniture and object clutter. Formally, given connected components $\mathcal{C} = \{C_k\}$ from the raw
 975 top-view occupancy, the VLM predicts a subset $\mathcal{C}_{wall} \subset \mathcal{C}$, and the wall map is constructed as

$$W = \bigcup_{C_k \in \mathcal{C}_{wall}} C_k.$$

976 The resulting room-level visual annotations include multi-height occupancy maps, object-layout
 977 maps, cleaned floor plans, and wall-structure maps. These representations provide global spatial
 978 priors for exploration, navigation target generation, and scene-level action annotation.

979 **Prompt for VLM-assisted wall filtering.** We use the following prompt to remove movable objects
 980 from the raw top-view occupancy map and retain structural room boundaries.

Wall-structure filtering prompt

System: You are a careful annotation engine for indoor scene structure. Your task is to identify persistent architectural structures in a top-view map. Keep walls, room boundaries, door frames, and other fixed structural elements. Remove movable objects such as tables, chairs, cabinets, shelves, sofas, lamps, and small objects.

User: You are given a top-view occupancy map with connected components and object labels. Each connected component has an ID. Please decide which components correspond to permanent wall or room-boundary structures.

Rules: 1. Keep components that form long, continuous room boundaries. 2. Keep components that separate rooms or define door openings. 3. Remove components that correspond to furniture, movable objects, or clutter. 4. If a component is uncertain, mark it as "uncertain" instead of guessing.

Return the result in JSON format with the following fields:
wall_component_ids, removed_object_component_ids, uncertain_component_ids, reasoning.

981

982 C Details of physics annotation pipeline

983 C.1 Details of the Unified Action Annotation Schema

984 In the main paper, we describe action annotation as a hierarchical candidate-bank schema. Here we
 985 provide additional details and concrete examples. The key idea is that action labels are not attached
 986 to an object as a single canonical action. Instead, they are organized as a **one-to-many mapping**
 987 from grounded visual anchors to compatible skills and then to multiple feasible action candidates:

visual anchor \rightarrow functional affordance \rightarrow compatible skills \rightarrow candidate bank \rightarrow action candidates.

988 This design explicitly preserves both **action diversity** and **functional consistency**. The visual-
 989 language annotations do not only identify where an object part or keypoint is; they also provide
 990 functional cues about how the region should be used. For example, a mug handle is not merely a
 991 graspable geometry: it is a functional handle that supports holding, lifting, and pouring. In contrast,
 992 the mug body may also be physically graspable, but such grasps are less functionally aligned for
 993 tasks that require using the mug as intended. Similarly, a cabinet handle supports both grasping and
 994 articulated opening, while the free space in front of the cabinet supports navigation and approach-pose
 995 annotation.

996 Given an asset \mathcal{A} , the visual annotation pipeline produces a set of grounded interaction anchors:

$$\mathcal{H}(\mathcal{A}) = \mathcal{P}(\mathcal{A}) \cup \mathcal{K}(\mathcal{A}) \cup \mathcal{R}^{\text{aff}}(\mathcal{A}) \cup \mathcal{R}^{\text{scene}}(\mathcal{A}),$$

997 where \mathcal{P} denotes part-level regions, \mathcal{K} denotes semantic keypoints, \mathcal{R}^{aff} denotes affordance regions,
 998 and $\mathcal{R}^{\text{scene}}$ denotes scene-level regions. Each anchor $h \in \mathcal{H}(\mathcal{A})$ is associated with a grounded
 999 functional affordance $\phi(h)$, such as grasp-for-holding, grasp-for-opening, articulation, hanging, or

1000 navigation approach. We use $\phi(h)$, together with language priors, geometric cues, and task constraints,
 1001 to infer a set of compatible skills $\mathcal{S}(h)$. For each anchor-skill pair, we maintain a candidate bank:

$$\mathcal{B}_{h,s} = \left\{ a_{h,s}^{(i)} \right\}_{i=1}^{N_{h,s}}, \quad s \in \mathcal{S}(h).$$

1002 Each candidate action is represented as

$$a_{h,s}^{(i)} = \left(s, o(h), h, \phi(h), x^{(i)}, \theta^{(i)}, \tau^{(i)}, v^{(i)}, d^{(i)} \right).$$

1003 The notation is summarized in Table 5.

1004 **Example: mug.** For a mug asset, the visual annotation pipeline may identify several anchors,
 1005 including the handle region, the mug body surface, and semantic keypoints on the handle or rim.
 1006 In this example, we focus on grasp-related skills. The mug handle is a **functional grasp anchor**:
 1007 grasping the handle preserves the intended use of the mug and supports downstream tasks such as
 1008 lifting, carrying, pouring, or placing. The mug body may also be physically graspable, but it is
 1009 not necessarily a functional grasp region for these long-horizon tasks, since grasping the body can
 1010 occlude the cup, interfere with pouring, or produce task-irrelevant contact. Therefore, the schema
 1011 distinguishes **functional grasp candidates** from merely **physically feasible grasp candidates**.

1012 For example, if $h = h_{\text{handle}}$ and $s = s_{\text{grasp}}$, then the candidate bank is

$$\mathcal{B}_{h_{\text{handle}},s_{\text{grasp}}} = \left\{ a_{h_{\text{handle}},s_{\text{grasp}}}^{(1)}, a_{h_{\text{handle}},s_{\text{grasp}}}^{(2)}, \dots, a_{h_{\text{handle}},s_{\text{grasp}}}^{(N_{h_{\text{handle}},s_{\text{grasp}}})} \right\}.$$

1013 These candidates may correspond to grasping the handle from different sides, using different gripper
 1014 orientations, or approaching with different pre-grasp offsets. Because the handle carries the affordance
 1015 $\phi(h_{\text{handle}}) = \text{grasp-for-use}$, these candidates are interpreted as functional grasp candidates rather
 1016 than merely stable grasps. In contrast, candidates from h_{body} may still be useful for simple pickup or
 1017 stabilization, but they are marked differently because they may not support downstream functional
 1018 tasks such as pouring. This allows downstream modules to select not only physically feasible grasps,
 1019 but also functionally appropriate grasps for long-horizon manipulation.

1020 **Example: cabinet.** For a cabinet or drawer asset, the visual annotation pipeline may identify
 1021 anchors such as the handle, door panel, drawer front, and reachable free-space region in front of the
 1022 cabinet. These anchors support different skills. A handle anchor may support both grasping and
 1023 articulation, while a scene-level free-space anchor supports navigation and approach-pose annotation.
 1024 The key point is that the handle is not simply a graspable part: it carries the functional affordance of
 1025 opening or closing the articulated structure.

1026 The cabinet example illustrates how functional grounding interacts with diversity. Diversity exists
 1027 **across anchors**: the handle and free-space region support different types of interaction. Diversity
 1028 also exists **within each anchor-skill pair**: the same handle-articulation pair may contain multiple
 1029 opening trajectories with different pull directions, waypoint sequences, and contact configurations.
 1030 At the same time, these candidates are function-conditioned by $\phi(h_{\text{handle}}) = \text{grasp-for-opening}$, so
 1031 a valid grasp or trajectory is not only collision-free and kinematically feasible, but also useful for
 1032 opening the cabinet or drawer. Therefore, the cabinet is not annotated with a single “open” action;
 1033 it is annotated with a structured bank of validated, function-conditioned grasp, articulation, and
 1034 navigation candidates.

1035 **Discussion.** This hierarchical representation provides a common interface for heterogeneous action
 1036 primitives while separating physical feasibility from functional suitability. Parallel-jaw grasping,
 1037 dexterous grasping, articulation, insertion, hanging, deformable manipulation, and navigation all
 1038 share the same high-level structure, while differing in the interpretation of $\phi(h)$, $x^{(i)}$, $\theta^{(i)}$, and $\tau^{(i)}$.
 1039 This allows downstream systems to query action annotations uniformly. For example, a manipulation
 1040 policy can request all functionally valid grasp candidates for a mug handle, a mobile manipulator
 1041 can request all feasible base poses for approaching a cabinet, and a data-collection module can
 1042 sample diverse candidates from $\mathcal{B}_{h,s}$ according to validation score, functional affordance, or diversity
 1043 descriptor.

Table 5: Symbols in the unified action annotation schema. The additional functional affordance term $\phi(h)$ distinguishes functionally meaningful actions from merely physically feasible ones.

| Symbol | Meaning | Examples |
|---|-----------------------|--|
| \mathcal{A} | Input asset | A single object such as a mug, or a room-scale scene containing a cabinet, table, and floor space. |
| $\mathcal{H}(\mathcal{A})$ | Grounded anchors | All actionable parts, keypoints, affordance regions, and scene-level regions detected on asset \mathcal{A} . |
| $\mathcal{P}(\mathcal{A})$ | Part regions | Mug handle, mug body, cabinet handle, cabinet door, drawer front. |
| $\mathcal{K}(\mathcal{A})$ | Semantic keypoints | Handle endpoint, drawer handle center, cabinet door corner, cloth corner. |
| $\mathcal{R}^{\text{aff}}(\mathcal{A})$ | Affordance regions | Functional handle region, graspable region, support edge, hook region, insertion socket. |
| $\mathcal{R}^{\text{scene}}(\mathcal{A})$ | Scene regions | Reachable free space in front of a cabinet, object-centric approach region, collision-free base-pose region. |
| h | Visual anchor | The handle region of a mug, the center of a cabinet handle, or a free-space region in front of a cabinet. |
| $\phi(h)$ | Functional affordance | For a mug handle: grasp-for-holding or grasp-for-pouring. For a mug body: physical grasp or stabilizing grasp. For a cabinet handle: grasp-for-opening. For free space: navigation approach. |
| $\mathcal{S}(h)$ | Compatible skills | For a mug handle: functional grasp or dexterous grasp. For a cabinet handle: grasp and articulation. For free space: navigation. |
| $\mathcal{B}_{h,s}$ | Candidate bank | All feasible function-conditioned grasps on a mug handle, or all feasible opening trajectories for a cabinet handle. |
| $N_{h,s}$ | Bank size | The number of validated grasp poses, dexterous contact sets, articulation trajectories, or base poses for skill s at anchor h . |
| $a_{h,s}^{(i)}$ | Action candidate | One grasp pose, one dexterous contact configuration, one articulation trajectory, or one navigation base pose. |
| s | Skill type | Grasp, dexterous grasp, articulation, insertion, hanging, deformable manipulation, or navigation. |
| $o(h)$ | Associated instance | The mug associated with a handle anchor, the cabinet door associated with a handle anchor, or the scene containing a navigation region. |
| $x^{(i)}$ | Concrete target | A contact point, surface patch, support edge, opening center, articulation handle point, or navigation base pose. |
| $\theta^{(i)}$ | Action parameters | 6D grasp pose, gripper width, contact set, hand configuration, insertion direction, hanging pose, articulation axis, or target base pose. |
| $\tau^{(i)}$ | Trajectory | Approach-contact-retreat trajectory for grasping, pull trajectory for articulation, or navigation approach path. |
| $v^{(i)}$ | Validation metadata | Collision status, inverse-kinematics feasibility, contact stability, task success score, path feasibility, trajectory statistics, and functional success. |
| $d^{(i)}$ | Diversity descriptor | Target location, approach direction, contact mode, trajectory family, grasp type, embodiment parameters, and perturbation seed. |

1044 C.2 Details of the General Physics-based Action Generation Pipeline

1045 C.2.1 Details of Candidate Target Generation

1046 In the main paper, candidate target generation is summarized as the first stage of the physics-based
 1047 action annotation pipeline. Here we provide additional details. This stage converts visual-language
 1048 annotations into concrete target anchors that seed the candidate banks. The overall process is:

skill s +language annotation+visual grounding \rightarrow candidate region $\mathcal{C}_s \rightarrow$ sampled target anchors \rightarrow filtered candidate bank.

Table 6: Example instantiations of the unified action annotation schema for a mug. The handle is treated as a functional grasp anchor, while the body surface mainly provides physically feasible but less functionally aligned grasp candidates.

| Symbol | Mug example | Interpretation |
|---|--------------------------|--|
| \mathcal{A} | Mug asset | The input object to be annotated, including its handle, body surface, rim, and opening geometry. |
| h_{handle} | Handle region | A part-level anchor from $\mathcal{P}(\mathcal{A})$ and an affordance anchor from $\mathcal{R}^{\text{aff}}(\mathcal{A})$. It provides the spatial support for functional grasp candidates. |
| $\phi(h_{\text{handle}})$ | Grasp-for-use | The functional affordance of the handle. It indicates that grasps on this anchor are intended for holding, lifting, carrying, or pouring the mug. |
| s_{grasp} | Parallel-jaw grasp | A compatible skill in $\mathcal{S}(h_{\text{handle}})$. This skill instantiates function-conditioned 6D gripper poses on the handle. |
| $\mathcal{B}_{h_{\text{handle}}, s_{\text{grasp}}}$ | Functional grasp bank | A set of feasible parallel-jaw grasp candidates on the mug handle, with different contact locations, approach directions, gripper orientations, and gripper widths. These candidates are both physically feasible and functionally aligned with using the mug. |
| $a_{h,s}^{(i)}$ | One handle grasp | One concrete functional grasp candidate. Here, $x^{(i)}$ is a contact point or patch on the handle, $\theta^{(i)}$ contains a 6D grasp pose and gripper width, $\tau^{(i)}$ stores an approach-contact-retreat trajectory, $v^{(i)}$ records IK, collision, and functional feasibility, and $d^{(i)}$ records approach direction and grasp mode. |
| s_{dex} | Dexterous grasp | Another compatible skill in $\mathcal{S}(h_{\text{handle}})$. The same functional handle anchor can also instantiate multi-finger contact configurations. |
| $\mathcal{B}_{h_{\text{handle}}, s_{\text{dex}}}$ | Functional dexgrasp bank | A set of feasible dexterous grasp candidates on the handle. Candidates differ in fingertip contacts, palm pose, joint configuration, and contact mode, while preserving the handle’s functional role. |
| h_{body} | Mug body surface | Another part-level visual anchor. It may be physically graspable, but it is generally less functionally aligned with using the mug for long-horizon tasks such as pouring or drinking. |
| $\phi(h_{\text{body}})$ | Physical grasp | The affordance of the body surface is treated as physical or stabilizing grasp rather than functional grasp-for-use. It can be useful for relocation or simple pickup, but may not preserve the intended mug functionality. |
| $\mathcal{B}_{h_{\text{body}}, s_{\text{dex}}}$ | Body grasp bank | A set of dexterous or stabilizing grasp candidates on the mug body surface, with different contact sets, hand configurations, and stability scores. These candidates may pass physical validation but receive different functional metadata from handle grasps. |

1049 The key point is that candidate targets are not sampled from the entire asset uniformly. Instead, they
1050 are sampled from skill-compatible and function-conditioned regions identified by the visual-language
1051 annotation pipeline. For example, a mug handle can be selected as a functional grasp region for
1052 holding or pouring, while the mug body can be selected as a physically graspable region for simple
1053 pickup or stabilization.

1054 **Candidate region construction.** Given an asset \mathcal{A} and a skill type s , we first construct a set of
1055 skill-compatible candidate regions:

$$\mathcal{C}_s(\mathcal{A}) = \Gamma_s(\mathcal{Y}^{\text{lang}}, \mathcal{Y}^{\text{vis}}, \mathcal{G}(\mathcal{A})),$$

1056 where $\mathcal{Y}^{\text{lang}}$ contains language-level semantics and functional priors, \mathcal{Y}^{vis} contains visual annotations
1057 such as part masks, keypoints, affordance regions, and scene maps, and $\mathcal{G}(\mathcal{A})$ denotes the asset
1058 geometry. Each candidate region $c \in \mathcal{C}_s(\mathcal{A})$ is associated with a functional affordance $\phi(c)$, such

1059 as grasp-for-use, physical grasp, grasp-for-opening, support-for-hanging, or navigation approach.
 1060 This functional affordance determines whether a region is merely physically usable or useful for a
 1061 downstream task.

1062 Different annotation sources provide candidate regions for different skills. Part masks are used to
 1063 identify object parts such as mug handles, mug bodies, cabinet handles, drawer fronts, and hooks.
 1064 Semantic keypoints are used for point-like anchors such as garment corners, handle centers, or
 1065 articulation keypoints. Affordance regions are used to identify functional regions such as handles,
 1066 openings, support edges, and hanging points. Room-level occupancy or BEV maps are used to
 1067 identify traversable free space and interaction-ready base-pose regions for navigation.

1068 **Sampling target anchors.** After constructing $\mathcal{C}_s(\mathcal{A})$, we sample concrete target anchors from each
 1069 candidate region. For object-centric skills, we first extract the subset of surface points or mesh faces
 1070 belonging to the candidate region. We then use farthest-point sampling (FPS) to obtain spatially
 1071 diverse target anchors:

$$p_1 \sim \mathcal{P}_c, \quad p_j = \arg \max_{p \in \mathcal{P}_c} \min_{1 \leq k < j} \|p - p_k\|_2,$$

1072 where \mathcal{P}_c is the set of surface points inside candidate region c . FPS avoids generating many redundant
 1073 anchors in the same local area and encourages coverage over the entire functional region. For each
 1074 sampled point, we estimate local geometric attributes such as surface normal, tangent directions,
 1075 curvature, local patch size, and local reference frame. These attributes are later used to instantiate
 1076 grasp poses, contact frames, insertion frames, or articulation waypoints.

1077 For keypoint-based skills, the annotated keypoints can directly serve as anchors. For example,
 1078 garment manipulation may use cloth corners or edge keypoints as grasp anchors. We can also sample
 1079 small perturbations around each keypoint to create nearby candidate anchors, then form bimanual
 1080 anchor pairs for skills such as stretching, folding, or lifting. These pairs are filtered by distance,
 1081 symmetry, visibility, and robot reachability.

1082 For scene-centric skills, we sample targets on occupancy or BEV maps. Let \mathcal{M}^{bev} denote the BEV
 1083 map and \mathcal{F} denote the set of traversable cells. We sample 2D positions from \mathcal{F} using grid sampling
 1084 or 2D FPS:

$$q_j = (x_j, y_j, \psi_j),$$

1085 where (x_j, y_j) is a reachable location and ψ_j is a base orientation. The orientation is chosen to face
 1086 the target object, affordance region, or manipulation anchor, possibly with multiple discrete approach
 1087 angles. These samples become navigation goals or interaction-ready base poses.

1088 **Geometry and embodiment filtering.** The sampled anchors are filtered before they are inserted
 1089 into the candidate bank. Filtering removes targets that are unlikely to produce physically feasible or
 1090 functionally meaningful actions. The filters depend on the skill.

1091 For grasping, we remove regions that are too small for the gripper, lack sufficient contact area, have
 1092 unsuitable surface normals, violate gripper-width limits, or cause immediate collision between the
 1093 gripper and the object. For example, if a part mask corresponds to a very thin or tiny decorative
 1094 detail, it may be visually detected but filtered out because the gripper cannot form a stable grasp
 1095 around it. For functional grasping, we further check whether the sampled target is consistent with the
 1096 intended affordance. A mug-handle grasp is retained as a functional grasp-for-use candidate, while a
 1097 mug-body grasp may be retained only as a physical or stabilizing grasp.

1098 For dexterous grasping, we filter contact sets that violate hand joint limits, produce self-collision, lack
 1099 sufficient fingertip support, or place contacts on semantically inappropriate regions. For articulation,
 1100 sampled handle points are filtered by whether they can support the required motion direction, maintain
 1101 contact during pulling or rotating, and avoid collision with nearby geometry. For insertion, sampled
 1102 openings are filtered by entry size, alignment feasibility, collision-free approach, and whether the
 1103 insertion axis is geometrically well defined. For hanging, support edges or hooks are filtered by edge
 1104 orientation, contact stability under gravity, and clearance around the hanging object.

1105 For navigation, BEV samples are filtered by traversability, clearance, path connectivity, and manipu-
 1106 lation reachability. A base pose must be collision-free, reachable from the scene layout, and close
 1107 enough for the robot to interact with the target object or part. We also remove base poses with poor
 1108 visibility or orientations that make the downstream manipulation target unreachable.

1109 **Forming candidate banks.** After sampling and filtering, each remaining target anchor becomes a
 1110 seed for action candidate generation. For a skill s and a candidate region c , the retained anchors are
 1111 denoted as

$$\mathcal{T}_{c,s} = \left\{ x_{c,s}^{(j)} \right\}_{j=1}^{M_{c,s}},$$

1112 where $x_{c,s}^{(j)}$ may be a contact point, local surface patch, support edge, opening frame, handle interaction
 1113 point, navigation goal, or interaction-ready base pose. These anchors are then used to instantiate the
 1114 candidate bank:

$$\mathcal{B}_{h,s} = \left\{ a_{h,s}^{(i)} \right\}_{i=1}^{N_{h,s}}.$$

1115 At this stage, the candidates mainly contain target-level information. Later stages fill in skill-specific
 1116 action parameters, trajectories, optimization results, and validation metadata. For example, a sampled
 1117 mug-handle anchor becomes several 6D grasp candidates with different approach directions and
 1118 gripper widths; a sampled cabinet-handle anchor becomes several grasp and articulation candidates;
 1119 and a sampled BEV base pose becomes a navigation candidate with an approach direction and
 1120 optional path.

1121 **Example: mug.** For a mug, language annotations identify the handle as a functional region for
 1122 holding, lifting, and pouring. The visual annotation pipeline grounds this cue into a handle part mask.
 1123 Candidate target generation then samples target anchors on the handle using FPS. These anchors seed
 1124 a functional grasp bank:

$$\mathcal{B}_{h_{\text{handle}},s_{\text{grasp}}}.$$

1125 In parallel, the mug body may be identified as a physically graspable surface. FPS on the body mask
 1126 can produce additional grasp anchors, but these are marked with a different functional affordance,
 1127 such as physical grasp or stabilizing grasp. Thus, both handle and body grasps can be physically
 1128 feasible, but only the handle grasps are labeled as functional grasp-for-use candidates for tasks such
 1129 as pouring.

1130 **Example: garment.** For a garment, semantic keypoints such as corners, sleeves, hems, or collar
 1131 points directly define candidate anchors. Instead of sampling densely over the entire cloth surface, the
 1132 pipeline samples around these keypoints and forms keypoint pairs for bimanual actions. For example,
 1133 two sleeve endpoints may seed a stretching trajectory, while two bottom corners may seed a folding
 1134 or lifting trajectory. Candidate pairs are filtered by visibility, distance, symmetry, and whether the
 1135 robot can reach both anchors simultaneously.

1136 **Example: cabinet.** For a cabinet, language annotations identify the handle as a functional region for
 1137 opening, and visual grounding provides the handle mask or handle keypoint. FPS samples interaction
 1138 anchors on the handle. The same handle region can seed both a grasp bank and an articulation bank:

$$\mathcal{B}_{h_{\text{handle}},s_{\text{grasp}}} \quad \text{and} \quad \mathcal{B}_{h_{\text{handle}},s_{\text{art}}}.$$

1139 Grasp candidates are filtered by gripper feasibility and collision, while articulation candidates are
 1140 filtered by whether the sampled handle point can support the required pull or rotation direction. For
 1141 a room-scale cabinet scene, the BEV map additionally provides free-space regions in front of the
 1142 cabinet. Grid sampling or 2D FPS generates candidate base poses, which are filtered by traversability,
 1143 clearance, and whether the robot can reach the handle from that pose.

1144 C.3 Details of Trajectory Generation

1145 After candidate target generation, each retained target anchor is converted into an initial trajectory
 1146 or waypoint sequence. This stage does not yet produce the final executable trajectory; instead, it
 1147 provides a template-level motion seed that is later refined by trajectory optimization and checked by
 1148 physics validation. Given a candidate action with anchor h , functional affordance $\phi(h)$, target $x^{(i)}$,
 1149 and skill-specific parameters $\theta^{(i)}$, we generate an initial trajectory

$$\tau_0^{(i)} = \mathcal{G}_s \left(h, \phi(h), x^{(i)}, \theta^{(i)}; \mathcal{A} \right),$$

1150 where \mathcal{G}_s is a primitive-specific trajectory generator. Although each skill has its own motion template,
 1151 we organize these templates by the source of their constraints, as summarized in Table 10.

1152 **Object-property-conditioned trajectories.** This family uses static object or garment properties
 1153 to generate motion templates. For parallel-jaw grasping, a candidate target provides a grasp pose
 1154 $T_{\text{grasp}} = [R_{\text{grasp}}, p_{\text{grasp}}]$ and an approach direction $\mathbf{a}_{\text{grasp}}$. We generate a short approach-contact-
 1155 retreat trajectory:

$$T_{\text{pre}} = [R_{\text{grasp}}, p_{\text{grasp}} - \delta_{\text{pre}}\mathbf{a}_{\text{grasp}}], \quad T_{\text{ret}} = [R_{\text{grasp}}, p_{\text{grasp}} + \delta_{\text{ret}}\mathbf{r}],$$

1156 where δ_{pre} is the pre-grasp offset, δ_{ret} is the retreat distance, and \mathbf{r} is a retreat direction such as the
 1157 reverse approach direction or an upward lifting direction. The resulting template is

$$\tau_0 = (T_{\text{pre}}, T_{\text{grasp}}, T_{\text{close}}, T_{\text{ret}}).$$

1158 For dexterous grasping, the same idea is extended to palm poses, fingertip contact sets, and hand joint
 1159 configurations. The pre-grasp pose is chosen to approach the object without early collision, and the
 1160 closure phase follows the candidate contact set.

1161 For garment manipulation, trajectories are generated from semantic keypoints. For an upper-body
 1162 garment, we use keypoints such as left/right sleeve endpoints, left/right bottom corners, and left/right
 1163 shoulder points:

$$\mathcal{K}_{\text{top}} = \{k_{\text{lslv}}, k_{\text{rslv}}, k_{\text{lb}}, k_{\text{rb}}, k_{\text{lsh}}, k_{\text{rsh}}\},$$

1164 where “slv” denotes sleeve, “b” denotes bottom, and “sh” denotes shoulder. A sleeve-fold trajectory
 1165 is generated by defining a fold line from the shoulder to the bottom corner on each side. For the left
 1166 sleeve, the fold line is

$$\ell_L = \ell(k_{\text{lsh}}, k_{\text{lb}}),$$

1167 and the sleeve target is obtained by reflecting the sleeve endpoint across this line:

$$k_{\text{lslv}}^* = \text{Ref}_{\ell_L}(k_{\text{lslv}}).$$

1168 Similarly, the right sleeve is folded using the line

$$\ell_R = \ell(k_{\text{rsh}}, k_{\text{rb}}), \quad k_{\text{rslv}}^* = \text{Ref}_{\ell_R}(k_{\text{rslv}}).$$

1169 The generated bimanual waypoint sequence lifts the sleeve endpoint, moves it toward the mirrored
 1170 target, and places it inside the garment body. After sleeve folding, we generate a bottom-up fold by
 1171 moving the bottom keypoints toward the shoulder line:

$$\ell_{\text{sh}} = \ell(k_{\text{lsh}}, k_{\text{rsh}}), \quad k_{\text{lb}}^* = \Pi_{\ell_{\text{sh}}}(k_{\text{lb}}), \quad k_{\text{rb}}^* = \Pi_{\ell_{\text{sh}}}(k_{\text{rb}}),$$

1172 where Π_{ℓ} denotes projection onto line ℓ . This produces a folding sequence in which the bottom hem
 1173 is lifted and folded toward the shoulder region.

1174 For pants, we use a symmetric folding rule. We first estimate the central axis of the pants from
 1175 waistband and leg keypoints. The two side or leg regions are folded toward the central axis so that the
 1176 legs overlap. Then the bottom or cuff keypoints are folded upward toward the waistband region. A
 1177 second bottom-up fold can be applied to compact the garment further. Thus, pants folding follows a
 1178 “side-fold then fold-again” template:

left/right side fold \rightarrow bottom-to-waist fold \rightarrow second compacting fold.

1179 For fling-like garment motions, we choose two keypoints, such as sleeve endpoints or bottom corners,
 1180 and generate a bimanual lift-and-stretch trajectory. The lift height, lateral stretch, and optional
 1181 shaking amplitude are scaled by the garment length

$$L_{\text{garment}} = \|k_1 - k_2\|_2,$$

1182 so that larger garments receive longer trajectories while remaining within bimanual reachability
 1183 constraints.

1184 Table 11 summarizes the garment templates.

1185 **Object-trajectory-conditioned trajectories.** This family is used for articulated objects whose
 1186 moving parts follow a joint trajectory. Examples include opening a cabinet door, pulling a drawer,
 1187 closing a lid, rotating a button, or pushing a button. The key constraint is that the end-effector should
 1188 follow the moving object part while maintaining the intended contact or grasp relation.

1189 Let $T_{\text{part}}(q)$ denote the pose of the moving part at joint state q , and let $T_{\text{eef}}^{\text{init}}$ denote the end-effector
 1190 pose after grasping or contacting the handle at the initial joint state q_{init} . For handle-based articulation,
 1191 we preserve the relative transform

$$T_{\text{rel}} = T_{\text{part}}(q_{\text{init}})^{-1}T_{\text{eef}}^{\text{init}}.$$

1192 For any later joint state q , the idealized end-effector pose is

$$T_{\text{eef}}(q) = T_{\text{part}}(q)T_{\text{rel}}.$$

1193 We generate a waypoint sequence by interpolating the joint state

$$q_j = q_{\text{init}} + j\Delta q, \quad j = 0, \dots, J,$$

1194 until a target joint state q_{goal} is reached. This formulation applies to both revolute joints, such as
 1195 doors and lids, and prismatic joints, such as drawers. The generated end-effector waypoints are later
 1196 checked by inverse kinematics, collision checking, and contact-maintenance validation.

1197 Different articulation actions use different joint schedules and contact assumptions, as shown in
 1198 Table 12.

1199 **Object-vector-conditioned trajectories.** This family is used when the action is defined by a
 1200 semantic or geometric direction rather than by an articulated part trajectory. Examples include
 1201 insertion, hanging, pouring, and placement. Each action has an active vector, a passive vector, or
 1202 both. The active vector belongs to the manipulated object, such as the axis of a peg or the pouring
 1203 direction of a container. The passive vector belongs to the target region, such as the axis of a socket,
 1204 the normal of a placement surface, or the support direction of a hook.

1205 For insertion, we align the active insertion axis \mathbf{u}_{act} with the passive insertion axis \mathbf{u}_{pas} :

$$R\mathbf{u}_{\text{act}} \approx -\mathbf{u}_{\text{pas}},$$

1206 where R is the end-effector or object orientation. Given an insertion target position p_{tar} and insertion
 1207 direction \mathbf{u}_{pas} , we generate an approach and execution trajectory:

$$p_{\text{pre}} = p_{\text{tar}} - \delta_{\text{pre}}\mathbf{u}_{\text{pas}}, \quad p(t) = p_{\text{pre}} + t\delta_{\text{ins}}\mathbf{u}_{\text{pas}}, \quad t \in [0, 1].$$

1208 For hanging, the end-effector approaches a support edge or hook along a clearance direction, places
 1209 the object into a stable contact configuration, and then lowers or releases the object along gravity.
 1210 For pouring, the trajectory rotates the grasped object so that its active pouring vector aligns with the
 1211 desired passive receiving direction or target container. For placement, the trajectory aligns the object
 1212 with the support normal, approaches the placement pose, releases, and retreats.

1213 Table 13 summarizes representative object-vector-conditioned templates.

1214 **Scene-trajectory-conditioned trajectories.** For scene-level skills, trajectory generation operates
 1215 on room-level occupancy or BEV maps. We first extract traversable free-space cells and sample
 1216 interaction-ready base poses around the target object, part, or affordance region. Each base pose is
 1217 represented as

$$q = (x, y, \psi),$$

1218 where (x, y) is a reachable location and ψ is the base orientation facing the manipulation target.
 1219 Given the current robot pose and a sampled target base pose, we compute a coarse global navigation
 1220 path using A* on the occupancy or BEV grid. This produces a sequence of grid-level waypoints
 1221 connecting free space to the interaction region.

1222 The A* path is used only as a geometric trajectory seed. It does not fully account for embodiment-
 1223 specific motion constraints, such as the turning radius of Ackermann platforms, differential-drive
 1224 dynamics, or holonomic constraints for quadrupeds and humanoids. These constraints are handled in
 1225 the subsequent trajectory optimization stage, where the A* path is locally refined by rollout-based
 1226 control such as DWB. The resulting scene-level trajectory provides a collision-free and interaction-
 1227 ready approach path to the manipulation target.

1228 **C.4 Details of Trajectory Optimization**

1229 After trajectory generation, each candidate action has an initial template-level trajectory $\tau_0^{(i)}$ and skill-
 1230 specific parameters $\theta_0^{(i)}$. These templates are intentionally simple: they are generated from geometric
 1231 rules, motion templates, A* paths, or keypoint-based heuristics. The trajectory optimization stage
 1232 refines them into asset-specific and embodiment-aware candidates before physics validation. Given a
 1233 candidate action

$$a_{h,s}^{(i)} = \left(s, o(h), h, \phi(h), x^{(i)}, \theta^{(i)}, \tau^{(i)}, v^{(i)}, d^{(i)} \right),$$

1234 we optimize the action parameters and trajectory as

$$\left(\theta_*^{(i)}, \tau_*^{(i)} \right) = \arg \min_{\theta, \tau} \mathcal{E}_s \left(\theta, \tau; h, \phi(h), x^{(i)}, \mathcal{A}, \mathcal{R} \right),$$

1235 where \mathcal{R} denotes the robot embodiment, including its kinematics, collision model, gripper or hand
 1236 model, and, for navigation, base dynamics. We use a skill-dependent objective:

$$\mathcal{E}_s = \lambda_{\text{func}} E_{\text{func}} + \lambda_{\text{task}} E_{\text{task}} + \lambda_{\text{geom}} E_{\text{geom}} + \lambda_{\text{contact}} E_{\text{contact}} + \lambda_{\text{coll}} E_{\text{coll}} + \lambda_{\text{kin}} E_{\text{kin}} + \lambda_{\text{smooth}} E_{\text{smooth}}.$$

1237 The terms are activated depending on the skill. E_{func} encourages consistency with the grounded
 1238 functional affordance $\phi(h)$, E_{task} tracks task-specific goals such as opening, insertion depth, or
 1239 navigation target, E_{geom} adapts waypoints to asset geometry, E_{contact} improves contact alignment or
 1240 contact maintenance, E_{coll} penalizes penetration and unsafe clearance, E_{kin} enforces embodiment-
 1241 specific kinematics or dynamics, and E_{smooth} regularizes motion. The optimized output is not yet
 1242 treated as a final annotation; it is passed to physics validation, where infeasible candidates are rejected.

1243 **Dexterous and contact-aware optimization.** For dexterous grasping, the initial candidate usually
 1244 specifies a target region, a set of possible contact points, and an initial hand or palm pose. However,
 1245 this template may not be physically plausible: fingertips may not align with the object surface, the
 1246 hand may penetrate the object, the joint configuration may violate limits, or the contact set may not
 1247 support the intended task. We therefore optimize the hand pose, finger joints, and contact assignment
 1248 jointly.

1249 Let q_{hand} denote hand joint angles, T_{hand} denote the hand root pose, and \mathcal{C} denote selected contact
 1250 links or fingertips. The contact point generated by forward kinematics for contact $c \in \mathcal{C}$ is denoted by

$$p_c = \text{FK}_c(T_{\text{hand}}, q_{\text{hand}}).$$

1251 Given a grounded functional region h with affordance $\phi(h)$, we use a region-consistency term

$$E_{\text{region}} = \sum_{c \in \mathcal{C}} \text{dist}(p_c, h)^2,$$

1252 where $\text{dist}(p_c, h)$ measures the distance from the contact point to the grounded part mask or af-
 1253 fordance region. This term prevents the optimizer from drifting to geometrically convenient but
 1254 functionally incorrect contacts. For example, if $\phi(h)$ is grasp-for-use on a mug handle, contacts are
 1255 encouraged to stay on the handle rather than moving to the mug body.

1256 To improve physical grasp quality, we add contact and wrench-compatibility terms. Let n_c be the
 1257 local surface normal at the matched object point and f_c be a feasible contact force under a friction
 1258 cone. A generic wrench-compatibility term can be written as

$$E_{\text{wrench}} = \sum_j \left\| w_j - \sum_{c \in \mathcal{C}} G_c f_{c,j} \right\|_2^2, \quad f_{c,j} \in \mathcal{K}_{\mu_c},$$

1259 where w_j are target wrench directions, G_c maps contact forces to object wrenches, and \mathcal{K}_{μ_c} is the
 1260 friction cone at contact c . This term encourages the selected contacts to support task-relevant external
 1261 wrenches. For a functional mug-handle grasp, the target wrench set can emphasize lifting or pouring
 1262 stability; for a cabinet-handle grasp, the target wrench set can emphasize pulling along the opening
 1263 direction.

1264 We also use distance, collision, and kinematic penalties:

$$E_{\text{contact}} = \sum_{c \in \mathcal{C}} \rho(d_{\mathcal{A}}(p_c)), \quad E_{\text{coll}} = \sum_b \max(0, \epsilon_{\text{safe}} - d_{\text{obs}}(b))^2,$$

1265 where d_A measures signed distance to the target object surface, $d_{\text{obs}}(b)$ measures clearance for hand
 1266 or arm body element b , and $\rho(\cdot)$ encourages near-surface contact without excessive penetration. Joint
 1267 limits and self-collision are included in E_{kin} . The optimized dexterous candidate stores the final hand
 1268 pose, joint configuration, contact set, pre-grasp pose, and approach trajectory in $\theta_*^{(i)}$ and $\tau_*^{(i)}$.

1269 **Pre-grasp and approach refinement.** For grasp-related skills, optimization also refines the pre-
 1270 grasp and approach motion. Given an optimized grasp pose T_{grasp} and approach direction \mathbf{a} , the
 1271 initial pre-grasp pose is typically placed as

$$T_{\text{pre}} = T_{\text{grasp}} \ominus \delta_{\text{pre}} \mathbf{a},$$

1272 where δ_{pre} is an offset along the reverse approach direction. This pre-grasp can fail when nearby
 1273 geometry blocks the approach path. We therefore optimize the approach offset, approach direction,
 1274 and intermediate waypoints to reduce collision and improve reachability:

$$E_{\text{approach}} = \sum_t E_{\text{coll}}(T_t) + \lambda_{\text{dir}} \sum_t \|\Delta p_t - \alpha_t \mathbf{a}\|_2^2 + \lambda_{\text{smooth}} \sum_t \|T_{t+1} \ominus T_t\|^2.$$

1275 This produces a collision-aware approach-contact-close-retreat trajectory. For dexterous hands, the
 1276 finger closure phase is also adjusted so that fingers contact the object progressively without early
 1277 penetration.

1278 **Adaptive geometry-aware waypoint optimization.** Some primitives require adapting template
 1279 waypoints to asset-specific geometry. This is especially important for deformable or articulated
 1280 objects, where a fixed motion template can fail across different shapes or scales. We represent the
 1281 generated template waypoints as

$$\tau_0 = \{p_1^0, \dots, p_T^0\},$$

1282 and optimize corrected waypoints

$$p_t^* = p_t^0 + \Delta p_t$$

1283 under task and geometry constraints:

$$\{\Delta p_t\}_{t=1}^T = \arg \min_{\{\Delta p_t\}} \sum_t \|\Delta p_t\|_2^2 + \lambda_{\text{geom}} E_{\text{geom}} + \lambda_{\text{reach}} E_{\text{reach}} + \lambda_{\text{smooth}} E_{\text{smooth}}.$$

1284 This keeps the optimized trajectory close to the template while adapting it to the current asset.

1285 For garment folding, we adapt fold targets according to keypoint distances. For an upper-body
 1286 garment, sleeve folding uses sleeve endpoints, shoulder keypoints, and bottom corners. The template
 1287 places the sleeve endpoint at a mirrored target inside the garment body. However, the exact place
 1288 point should depend on sleeve length and body width. We therefore scale the inward displacement by
 1289 garment geometry:

$$\delta_{\text{sleeve}} = \text{clip}(\alpha_{\text{sleeve}} L_{\text{sleeve}}, \delta_{\text{min}}, \delta_{\text{max}}),$$

1290 where L_{sleeve} is the distance between the shoulder keypoint and sleeve endpoint. The sleeve place
 1291 point is shifted inward by δ_{sleeve} along the fold direction. Similarly, for bottom-up folding, the target
 1292 position is adapted using the hem-to-shoulder distance:

$$\delta_{\text{bottom}} = \text{clip}(\alpha_{\text{bottom}} L_{\text{hem} \rightarrow \text{shoulder}}, \delta_{\text{min}}, \delta_{\text{max}}).$$

1293 This prevents short garments from being over-folded and long garments from being under-folded.

1294 For pants, the generated template first folds the two sides or legs toward the centerline and then folds
 1295 the bottom region upward toward the waistband. Optimization adjusts the side-fold and bottom-fold
 1296 targets according to leg length, waistband width, and bimanual reachability. If the two grasp points
 1297 are too far apart for the embodiment, the optimizer moves them to nearby keypoint perturbations or
 1298 selects a shorter folding displacement while preserving the same folding intention.

1299 For fling-like garment motions, the template chooses two keypoints and generates a bimanual lift-and-
 1300 stretch trajectory. The lift height and lateral stretch are scaled by the distance between the selected
 1301 keypoints:

$$H_{\text{lift}} = \text{clip}(\alpha_H \|k_1 - k_2\|_2, H_{\text{min}}, H_{\text{max}}), \quad D_{\text{stretch}} = \text{clip}(\alpha_D \|k_1 - k_2\|_2, D_{\text{min}}, D_{\text{max}}).$$

1302 The optimizer then adjusts these values to satisfy bimanual reachability and avoid self-collision
 1303 between the two arms.

1304 **Articulation trajectory refinement.** For articulated objects, trajectory generation produces way-
 1305 points by following the motion of a movable part, such as a cabinet door, drawer, lid, or button. The
 1306 initial trajectory assumes that the robot maintains an intended contact or grasp relation to the moving
 1307 part. However, a direct template may fail because of robot reachability limits, collision with the
 1308 object frame, or loss of contact. We therefore optimize the end-effector waypoints, joint schedule,
 1309 and contact pose.

1310 Let q_t denote the object joint state at waypoint t , and let $T_{\text{eff},t}$ be the corresponding end-effector
 1311 pose. We use a tracking term that encourages the end-effector to follow the moving part:

$$E_{\text{art-track}} = \sum_t \left\| T_{\text{eff},t} \ominus \widehat{T}_{\text{eff}}(q_t) \right\|^2,$$

1312 where $\widehat{T}_{\text{eff}}(q_t)$ is the idealized end-effector pose computed from the articulated part motion and the
 1313 initial grasp or contact relation. We also use a contact-maintenance term:

$$E_{\text{maintain}} = \sum_t \text{dist}(p_{\text{contact},t}, h)^2 + \lambda_n \sum_t (1 - n_{\text{contact},t}^\top n_{\text{part},t}),$$

1314 which keeps the contact point on the handle or movable part and encourages consistent contact
 1315 normals.

1316 For handle-based opening and closing, the optimizer adjusts the waypoint spacing and pull direction
 1317 so that the robot can maintain contact while the door or drawer moves. For push-based closing,
 1318 the trajectory is optimized to maintain safe pushing contact without requiring a fixed grasp. For
 1319 buttons, the optimization bounds the push displacement by the button travel range and aligns the push
 1320 direction with the button normal. For rotating buttons or lids, angular waypoints are refined to satisfy
 1321 wrist reachability and avoid collisions.

1322 **Object-vector-conditioned trajectory refinement.** Vector-conditioned skills are defined by se-
 1323 mantic or geometric directions, such as insertion axes, hanging directions, pouring directions, and
 1324 placement normals. The template trajectory aligns the object or end-effector with the vector and
 1325 moves along it. Optimization improves alignment, clearance, and task-specific feasibility.

1326 For insertion, let \mathbf{u}_{act} be the active insertion axis of the manipulated object and \mathbf{u}_{pas} be the passive
 1327 insertion axis of the target opening or socket. The alignment cost is

$$E_{\text{align}} = 1 - (R\mathbf{u}_{\text{act}})^\top (-\mathbf{u}_{\text{pas}}),$$

1328 where R is the object or end-effector orientation. The insertion trajectory also includes an approach-
 1329 clearance term and an insertion-depth term:

$$E_{\text{insert}} = \lambda_{\text{align}} E_{\text{align}} + \lambda_{\text{depth}} (z_{\text{target}} - z_{\text{inserted}})^2 + \lambda_{\text{coll}} E_{\text{coll}}.$$

1330 This refines the alignment pose, entry pose, and insertion depth.

1331 For hanging, optimization adjusts the approach pose so that the object clears the hook or support edge,
 1332 establishes contact, and lowers into a gravity-stable configuration. The stability term encourages the
 1333 object center of mass to lie below or behind the support contact after release:

$$E_{\text{hang}} = E_{\text{align}} + \lambda_{\text{com}} \max(0, \Delta_{\text{unstable}})^2 + \lambda_{\text{clear}} E_{\text{clearance}}.$$

1334 For pouring, the optimizer refines the rotation trajectory around the grasp frame so that the pouring
 1335 vector points toward the receiving region while avoiding excessive wrist rotation or collision. For
 1336 placement, the optimizer aligns the object with the support normal and adjusts the release height and
 1337 retreat direction.

1338 **Scene-level trajectory optimization with DWB.** For navigation and approach-pose annotation,
 1339 trajectory generation first produces a coarse global path using A* on an occupancy or BEV map.
 1340 This A* path is geometrically valid on the grid, but it may not respect the robot’s embodiment
 1341 dynamics. For example, an Ackermann platform has a minimum turning radius, a differential-drive
 1342 base has coupled forward and angular velocities, while a quadruped or humanoid base can often
 1343 be approximated as holonomic in the navigation layer. We therefore refine the A* path using a
 1344 DWB-style local rollout optimizer.

1345 Let the robot base state be

$$q_t = (x_t, y_t, \psi_t),$$

1346 and let u_t be a velocity command. The admissible control set depends on the embodiment:

$$u_t \in \mathcal{U}_{\mathcal{R}}.$$

1347 For a differential-drive base,

$$u_t = (v_t, \omega_t),$$

1348 and the rollout model is

$$x_{t+1} = x_t + \Delta t v_t \cos \psi_t, \quad y_{t+1} = y_t + \Delta t v_t \sin \psi_t, \quad \psi_{t+1} = \psi_t + \Delta t \omega_t.$$

1349 For an Ackermann platform,

$$u_t = (v_t, \delta_t), \quad \dot{\psi}_t = \frac{v_t}{L} \tan \delta_t, \quad |\delta_t| \leq \delta_{\max},$$

1350 where L is the wheelbase and δ_t is the steering angle. This naturally enforces a turning-radius
1351 constraint. For holonomic bases, such as quadruped or humanoid navigation abstractions, we use

$$u_t = (v_{x,t}, v_{y,t}, \omega_t),$$

1352 which allows lateral motion in addition to rotation.

1353 DWB samples short-horizon velocity commands and rolls them out under the corresponding dynamics.
1354 Each rollout is scored by

$$J_{\text{DWB}} = w_{\text{path}} J_{\text{path}} + w_{\text{goal}} J_{\text{goal}} + w_{\text{obs}} J_{\text{obs}} + w_{\text{clear}} J_{\text{clear}} + w_{\text{smooth}} J_{\text{smooth}} + w_{\text{reach}} J_{\text{reach}}.$$

1355 Here, J_{path} measures deviation from the A* path, J_{goal} measures distance to the target base pose,
1356 J_{obs} penalizes collision with obstacles, J_{clear} encourages clearance, J_{smooth} penalizes abrupt velocity
1357 changes, and J_{reach} evaluates whether the final base pose can support manipulation of the target
1358 object or part. The best rollout becomes the optimized local trajectory segment. Repeating this
1359 process along the path converts the A* seed into an embodiment-feasible approach trajectory.

1360 **Manipulation reachability during navigation optimization.** For scene-level trajectories, reaching
1361 a base pose is not sufficient; the pose must also support the downstream manipulation target. We
1362 therefore add a reachability score to the DWB objective. Given a target manipulation anchor h and a
1363 candidate base pose q , we estimate whether the robot arm or body can reach the anchor:

$$J_{\text{reach}}(q, h) = \min_{\xi \in \Xi(q)} \|\text{FK}(\xi) - x_h\|^2 + \lambda_{\text{ik}} \mathbf{1}_{\text{IK fail}},$$

1364 where $\Xi(q)$ is the set of feasible arm configurations at base pose q , and x_h is the target position or
1365 pose induced by anchor h . In practice, this can be approximated using reachability maps, IK queries,
1366 or distance-to-workspace heuristics. This term favors base poses from which the robot can both see
1367 and manipulate the target part.

1368 **Candidate scoring and metadata.** After optimization, each candidate receives an optimization
1369 score and metadata before physics validation. We store both the optimized action and the individual
1370 objective terms:

$$v_{\text{opt}}^{(i)} = \{E_{\text{func}}, E_{\text{task}}, E_{\text{geom}}, E_{\text{contact}}, E_{\text{coll}}, E_{\text{kin}}, E_{\text{smooth}}, E_{\text{total}}\}.$$

1371 These values are later combined with simulation-based validation results. The diversity descriptor
1372 $d^{(i)}$ is also updated after optimization to record the final target location, approach direction, contact
1373 mode, trajectory family, base-pose class, or embodiment-specific command profile.

1374 **Discussion.** The optimization stage is deliberately separated from physics validation. Optimization
1375 improves candidate quality under differentiable or efficiently computable objectives, while validation
1376 performs stricter checks such as simulation rollout, collision verification, IK success, traversability,
1377 contact stability, and task-specific success. This separation allows AnnotateAnything to scale to
1378 many candidates per anchor-skill pair while still retaining only physically feasible and functionally
1379 meaningful action annotations.

1380 **C.5 Details of Physics Validation**

1381 Physics validation is the final filtering stage before an optimized candidate is stored as an action an-
 1382 notation. Although candidate generation and trajectory optimization produce geometrically plausible
 1383 actions, many candidates still fail when executed in simulation due to unstable contact, insufficient
 1384 force closure, object penetration, articulation blockage, unreachable arm motion, or deformable-object
 1385 entanglement. Therefore, validation is a critical step for improving downstream rollout success.

1386 Given an optimized candidate

$$a_{h,s}^{(i)} = \left(s, o(h), h, \phi(h), x^{(i)}, \theta_*^{(i)}, \tau_*^{(i)}, v^{(i)}, d^{(i)} \right),$$

1387 we run physics rollouts under a validation embodiment e_{val} and a perturbation setting Ξ :

$$v^{(i)} = \mathcal{V}_s \left(a_{h,s}^{(i)}, \mathcal{A}, e_{\text{val}}, \Xi \right).$$

1388 The validation result $v^{(i)}$ records success flags, collision statistics, contact stability, task progress,
 1389 robustness scores, rollout length, disturbance settings, and the embodiment used for validation.
 1390 Candidates that fail the corresponding skill-specific checks are rejected.

1391 **Floating embodiment validation.** For most object-centric skills, we use a **floating** embodiment,
 1392 such as a floating parallel gripper or floating dexterous hand. The floating embodiment directly
 1393 executes the optimized object-centric trajectory without attaching the gripper or hand to a specific
 1394 robot arm or mobile base. This setting is used for skills where the main question is whether the
 1395 local object interaction is physically meaningful, such as parallel-jaw grasping, dexterous grasping,
 1396 insertion, hanging, placement, pouring, and many handle-based articulation candidates.

1397 This design decouples object-level physical feasibility from a particular downstream robot config-
 1398 uration. In downstream data collection, robot-object relative poses may be randomized, and the
 1399 same asset-level annotation can be consumed by different robot embodiments. Enforcing full-arm
 1400 inverse kinematics during annotation would prematurely reject many locally valid interactions simply
 1401 because one particular arm pose cannot reach them. Floating validation instead focuses on contact,
 1402 penetration, local collision, object stability, task progress, and whether the action preserves the
 1403 intended functional affordance $\phi(h)$.

1404 For a parallel-jaw grasp, the floating gripper follows the pre-grasp trajectory, closes at the grasp pose,
 1405 and then runs a stability test. For a dexterous grasp, the floating hand moves to the pre-grasp pose,
 1406 closes according to the optimized joint configuration or contact schedule, and checks whether the
 1407 object is actually constrained by the hand. For articulation, the floating gripper or hand grasps or
 1408 contacts the handle and follows the generated opening, closing, rotating, or pushing trajectory. For
 1409 vector-conditioned skills such as insertion or hanging, the floating end-effector follows the aligned
 1410 approach and execution trajectory and verifies whether the object reaches a stable or task-complete
 1411 state.

1412 **Non-floating or full-robot validation.** Some skills are strongly embodiment-dependent and cannot
 1413 be reliably validated with a floating gripper or hand. For these skills, we use a **non-floating** validation
 1414 setting, including the full manipulator, arm, gripper or hand, and, when needed, the mobile base.
 1415 This setting is used when IK feasibility, arm reachability, self-collision, singularities, or robot-body
 1416 interaction directly determine task success.

1417 We use full-manipulator validation for tasks such as garment pushing, washing, retrieval from
 1418 entangled configurations, and other long-horizon deformable-object interactions. These tasks can
 1419 fail even when the local hand trajectory looks reasonable, because the arm may pass through an IK
 1420 singularity, become unreachable, collide with the environment, or entangle with the garment. In
 1421 such cases, floating validation would overestimate feasibility. Full-robot validation therefore checks
 1422 arm IK, joint limits, self-collision, arm-object collision, singular configurations, and task-specific
 1423 deformable outcomes.

1424 For navigation and approach-pose annotations, we use a base-aware non-floating setting. The A^*
 1425 path generated in the trajectory generation stage is checked and locally refined under the robot’s
 1426 base model, for example through DWB-style rollout. This allows the validation to account for
 1427 embodiment-specific constraints such as differential-drive velocity limits, Ackermann turning radius,
 1428 or holonomic base motion.

1429 **Disturbance and gravity robustness.** For grasp and dexterous grasp candidates, nominal execution
1430 is often insufficient. A grasp can appear successful under one gravity direction but fail as soon as
1431 the object is perturbed or the gravity direction changes. To reject such fragile candidates, we apply
1432 randomized disturbances and gravity robustness tests during validation.

1433 For each grasp or dexterous grasp candidate, we sample $K = 8$ random gravity directions:

$$\mathbf{g}_k = 1.5g_0\mathbf{u}_k, \quad \mathbf{u}_k \sim \mathbb{S}^2, \quad k = 1, \dots, 8,$$

1434 where g_0 is the nominal gravity magnitude and \mathbf{u}_k is a random unit direction. The gravity magnitude
1435 is scaled by 1.5 to create a stricter robustness test. We also apply small randomized disturbances,
1436 such as object pose jitter, hand pose jitter, or external perturbation forces, depending on the skill and
1437 simulator setting. A candidate is accepted only if it remains stable under the required robustness tests
1438 or exceeds the skill-specific pass threshold.

1439 This gravity randomization is especially useful for detecting grasps with insufficient contact. For
1440 example, a parallel-jaw grasp may lift the object under the default gravity direction but drop it when
1441 gravity is rotated, indicating that the contact patch is too small or the object is only supported by
1442 incidental contact. Similarly, a dexterous grasp may visually appear to wrap around the object, but
1443 the object may not actually be constrained by the fingers; randomized gravity reveals such failures.

1444 **Skill-specific validation criteria.** Each skill family uses task-specific success criteria in addition to
1445 generic collision and stability checks.

1446 For parallel-jaw grasping, we check whether the gripper can approach without collision, close on
1447 the object, maintain sufficient contact, and keep the object stable under disturbance and randomized
1448 gravity. The object should not slip, fall, or rotate beyond the allowed tolerance during the stability
1449 rollout.

1450 For dexterous grasping, we check hand joint validity, self-collision, hand-object penetration, fingertip
1451 or palm contact, and object stability. We also check whether the optimized grasp actually constrains
1452 the object rather than merely placing fingers near the object. This distinction is important because
1453 dexterous grasp candidates can look visually plausible while providing little or no effective contact.

1454 For articulation, we check whether the end-effector maintains contact with the handle or moving part,
1455 whether the object joint progresses toward the intended open, close, rotate, or push state, and whether
1456 the gripper or hand becomes stuck in the geometry. For dexterous articulation, we additionally check
1457 hand joint limits and self-collision, since fingers can become trapped in narrow handles or deformed
1458 by contact with the moving part.

1459 For insertion, hanging, placement, and pouring, we check alignment, clearance, stability after release,
1460 and task progress. For example, insertion candidates should reach sufficient insertion depth without
1461 collision; hanging candidates should remain stable under gravity after release; placement candidates
1462 should settle on the support surface; and pouring candidates should maintain the intended orientation
1463 and avoid collision.

1464 For garment and deformable-object skills, we check full-arm reachability, IK stability, collision
1465 with the cloth or environment, and whether the resulting cloth state satisfies the intended geometric
1466 outcome. We reject trajectories that pass through singularities, cause the arm to move erratically,
1467 wrap the garment around the arm, or lead to severe entanglement.

1468 For navigation, we check that the path is traversable, collision-free, connected to the target base pose,
1469 and compatible with the robot’s base model. We also check manipulation reachability at the final
1470 base pose, since a navigation target is useful only if the robot can interact with the annotated object
1471 or part from that pose.

1472 **Parallel validation.** Physics validation is computationally expensive because each asset can contain
1473 many anchors, each anchor can support multiple skills, and each anchor-skill pair can store many
1474 candidates. Moreover, the pass rate can be low, especially for dexterous grasping, articulation,
1475 and deformable-object tasks. We therefore run validation in parallel across assets, skills, anchors,
1476 candidates, and perturbation trials.

1477 The validation jobs are grouped by skill and embodiment so that candidates with similar simulation
1478 settings can be batched together. For example, floating-gripper grasp candidates can be evaluated
1479 in one group, floating-dexterous-hand candidates in another group, and full-manipulator garment

1480 trajectories in a separate group. Within each group, multiple candidates and multiple robustness trials
 1481 can be executed concurrently. We also use early termination: if a candidate penetrates the object,
 1482 loses contact, drops the object, fails to progress, or violates a joint limit, the rollout is stopped and
 1483 marked as failed. Parallel validation is essential for scalability; without it, the number of simulation
 1484 rollouts would become prohibitive.

1485 **Common failure cases.** The validation stage rejects many candidates that appear plausible from
 1486 geometry alone. Table 22 summarizes common failure modes observed across different skills.

1487 **Discussion.** The separation between floating and non-floating validation is important for scalability
 1488 and generality. Floating validation keeps the annotation object-centric and reusable across downstream
 1489 embodiments, while full-robot validation is applied only when embodiment feasibility is intrinsic to
 1490 the skill. Robustness tests, especially randomized gravity directions and disturbance rollouts, prevent
 1491 the dataset from retaining fragile candidates that only succeed under a single idealized condition.
 1492 Finally, parallel validation makes it practical to evaluate large candidate banks and retain high-quality,
 1493 physically feasible, and functionally meaningful action labels.

1494 C.6 Details of Physics-aware Augmentation

1495 After physics validation, each accepted candidate is physically feasible under its original target, pose,
 1496 trajectory, and validation setting. However, storing only the validated seed candidates would still
 1497 limit action diversity. We therefore apply **physics-aware augmentation** to expand each candidate
 1498 bank while preserving physical feasibility and functional consistency. The goal is to produce multiple
 1499 valid variations of an action around the same functional anchor, rather than changing the semantic
 1500 meaning of the action.

1501 Given a validated candidate

$$a_{h,s}^{(i)} = \left(s, o(h), h, \phi(h), x^{(i)}, \theta^{(i)}, \tau^{(i)}, v^{(i)}, d^{(i)} \right),$$

1502 we generate an augmented set

$$\mathcal{U}\left(a_{h,s}^{(i)}\right) = \left\{ \tilde{a}_{h,s}^{(i,j)} \right\}_{j=1}^{M_i},$$

1503 where each augmented candidate $\tilde{a}_{h,s}^{(i,j)}$ preserves the same skill s , visual anchor h , and functional
 1504 affordance $\phi(h)$, but perturbs the concrete target, action parameters, or trajectory. The augmented
 1505 candidate bank becomes

$$\tilde{\mathcal{B}}_{h,s} = \mathcal{B}_{h,s} \cup \bigcup_{a_{h,s}^{(i)} \in \mathcal{B}_{h,s}} \mathcal{U}\left(a_{h,s}^{(i)}\right).$$

1506 We use two complementary augmentation mechanisms: **local perturbation augmentation** and
 1507 **symmetry-aware augmentation**.

1508 **Local perturbation augmentation.** Local perturbation samples small variations around a validated
 1509 anchor target or trajectory while keeping the action inside the same grounded region. This is useful
 1510 because many skills admit a local family of feasible actions. For example, a mug handle can be
 1511 grasped at slightly different contact centers and approach rotations; a garment sleeve fold can place
 1512 the sleeve endpoint within a small region around the nominal mirrored target; and a cabinet handle
 1513 can be pulled from slightly different handle points or approach directions.

1514 For a target $x^{(i)}$, we sample a perturbed target

$$\tilde{x}^{(i,j)} = x^{(i)} + \Delta x^{(j)}, \quad \Delta x^{(j)} \sim \mathcal{D}_x(h, s),$$

1515 where $\mathcal{D}_x(h, s)$ is a bounded perturbation distribution defined by the anchor type and skill. For
 1516 surface-based anchors, the perturbation is restricted to the part mask or affordance region. For
 1517 keypoint-based anchors, the perturbation is restricted to a keypoint confidence region. For scene-level
 1518 anchors, the perturbation is restricted to traversable free space or a reachable base-pose region.

1519 For pose-based actions, we also perturb orientation and skill parameters:

$$\tilde{\theta}^{(i,j)} = \theta^{(i)} \oplus \Delta \theta^{(j)},$$

1520 where $\Delta\theta^{(j)}$ may include in-plane grasp rotation, approach offset, gripper width adjustment, dexterous
 1521 hand contact reassignment, insertion-axis offset, pulling-direction perturbation, or base-orientation
 1522 perturbation. For trajectory-based actions, we perturb waypoint locations or timing:

$$\tilde{\tau}^{(i,j)} = \left\{ T_t \oplus \Delta T_t^{(j)} \right\}_{t=1}^T.$$

1523 These perturbations are bounded to preserve the original task intent. For example, a cabinet-handle
 1524 articulation candidate may perturb the handle contact point and pulling direction, but it must still
 1525 move along the valid opening direction. Similarly, a functional mug-handle grasp may perturb the
 1526 contact center and gripper yaw, but it should remain on the handle and preserve the grasp-for-use
 1527 affordance.

1528 **Garment-specific local augmentation.** Garment actions benefit from local perturbation because
 1529 keypoint annotations and deformable-object outcomes naturally contain uncertainty. For fling, we
 1530 jitter the selected grasp keypoints within local confidence regions around sleeve endpoints, bottom
 1531 corners, or hem keypoints. The lift height, lateral stretch, and release location can also be randomly
 1532 perturbed within a bounded range scaled by garment size. For folding, we perturb the fold line
 1533 and place target. For example, in upper-body garment folding, the sleeve endpoint is first mirrored
 1534 inward using the shoulder–bottom fold line; augmentation then samples nearby place points around
 1535 this mirrored target. For bottom-up folding, the bottom-corner place points are jittered around the
 1536 shoulder or upper-body target region. For pants, the side-fold target and bottom-up fold target are
 1537 perturbed around the centerline and waistband region. These variations increase trajectory diversity
 1538 while preserving the intended fold structure.

1539 **Grasp and dexterous grasp augmentation.** For parallel-jaw grasping, we augment a validated
 1540 grasp by perturbing the contact center, approach offset, gripper width, and in-plane rotation around
 1541 the grasp axis:

$$\tilde{T}_{\text{grasp}} = T_{\text{grasp}} \oplus (\Delta p, \Delta R),$$

1542 where Δp is constrained to remain within the graspable part or affordance region, and ΔR is bounded
 1543 to avoid changing the grasp family. For dexterous grasping, we perturb palm pose, fingertip contact
 1544 seeds, finger closure timing, and small joint offsets. These perturbations must preserve contact
 1545 coverage, joint-limit feasibility, self-collision constraints, and functional-region consistency. For
 1546 instance, dexterous mug-handle grasps remain on the handle, while body grasps remain labeled as
 1547 physical or stabilizing grasps rather than functional grasp-for-use.

1548 **Articulation augmentation.** For articulation, local perturbation is applied to the handle contact
 1549 point, grasp orientation, pulling direction, waypoint spacing, and target joint range. This is important
 1550 because a single handle can often be pulled from multiple nearby points and directions. However, the
 1551 perturbation must preserve the articulation affordance: opening trajectories should still follow the
 1552 joint motion, closing trajectories should move toward the closed state, and push-button trajectories
 1553 should remain aligned with the button normal. Perturbed articulation candidates are checked for
 1554 collision, contact maintenance, joint progress, and whether the gripper or dexterous hand becomes
 1555 stuck in the handle geometry.

1556 **Navigation and approach-pose augmentation.** For scene-level annotations, we perturb interaction-
 1557 ready base poses and approach directions within reachable free-space regions. Given a base pose

$$q = (x, y, \psi),$$

1558 we sample

$$\tilde{q} = (x + \Delta x, y + \Delta y, \psi + \Delta \psi),$$

1559 where $(x + \Delta x, y + \Delta y)$ must remain in traversable free space and $\psi + \Delta \psi$ should keep the robot
 1560 oriented toward the target object or affordance region. The augmented base pose is rechecked for
 1561 clearance, path connectivity, DWB feasibility, and manipulation reachability. This produces multiple
 1562 approach poses around the same object, enabling downstream policies to sample different viewpoints
 1563 and approach directions.

1564 **Symmetry-aware augmentation.** Local perturbation only explores nearby variations. To further
 1565 expand diversity, we use symmetry-aware augmentation when visual-language annotations indicate
 1566 that an object or part has a valid symmetry. Let $\mathcal{G}_{\text{sym}}(o)$ denote the annotated symmetry group of
 1567 object or part o . For a validated candidate, we generate transformed candidates

$$\tilde{a}_{h,s}^{(i,g)} = g \cdot a_{h,s}^{(i)}, \quad g \in \mathcal{G}_{\text{sym}}(o),$$

1568 where the symmetry transform g is applied consistently to the target $x^{(i)}$, pose parameters $\theta^{(i)}$, and
 1569 trajectory waypoints $\tau^{(i)}$. For pose-based actions,

$$\tilde{T} = GT,$$

1570 where G is the rigid transformation induced by the symmetry. For trajectory-based actions,

$$\tilde{\tau} = \{GT_1, GT_2, \dots, GT_T\}.$$

1571 For vector-conditioned actions, the relevant vectors are also transformed:

$$\tilde{\mathbf{u}} = R_G \mathbf{u},$$

1572 where R_G is the rotational component of G .

1573 For example, a bottle or can may have approximate rotational symmetry around its vertical axis.
 1574 If one grasp pose is validated, rotating the grasp around the symmetry axis can produce additional
 1575 valid grasps. A round knob may support rotationally symmetric grasp or turning candidates. A
 1576 drawer with repeated identical handles may support translated or mirrored handle grasps if the
 1577 visual-language annotation marks the handles as functionally equivalent. Garments may also contain
 1578 bilateral structure, such as left and right sleeves, allowing sleeve-fold templates to be mirrored
 1579 across the garment centerline. However, symmetry augmentation is applied only when it preserves
 1580 the functional affordance. A mug with a handle is not treated as fully rotationally symmetric for
 1581 functional grasping, because rotating a handle grasp around the mug body would move it away from
 1582 the handle and destroy the grasp-for-use affordance.

1583 **Feasibility checks for augmented candidates.** Augmented candidates are not automatically ac-
 1584 cepted. Each candidate is rechecked using lightweight geometry and physics constraints, and high-risk
 1585 skills can be sent back to full physics validation. The checks include:

$$\mathbb{I}_{\text{accept}} = \mathbb{I}[C_{\text{region}} \wedge C_{\text{func}} \wedge C_{\text{coll}} \wedge C_{\text{kin}} \wedge C_{\text{task}}],$$

1586 where C_{region} checks that the perturbed target remains inside the correct part mask, affordance region,
 1587 keypoint confidence region, or traversable scene region; C_{func} checks that the functional affordance
 1588 $\phi(h)$ is preserved; C_{coll} checks collision and clearance; C_{kin} checks embodiment constraints such
 1589 as gripper width, hand joint limits, or base dynamics; and C_{task} checks task-specific validity, such
 1590 as grasp stability, articulation direction, fold geometry, insertion alignment, hanging support, or
 1591 navigation reachability.

1592 For inexpensive skills, such as grasp or navigation pose augmentation, we can generate many
 1593 augmented candidates and filter them by geometry, collision, and short rollout checks. For expensive
 1594 skills, such as full-manipulator garment trajectories, we use fewer augmentations and run stricter
 1595 validation because small perturbations can produce IK singularities, arm-cloth entanglement, or large
 1596 deviations in deformable-object state.

1597 **Examples.** For a bottle, the visual-language annotation may mark the object as approximately
 1598 rotationally symmetric around the vertical axis. A validated side grasp can then be rotated by multiple
 1599 angles around this axis to produce additional grasp candidates:

$$G_\alpha = \text{Rot}_z(\alpha), \quad \alpha \in \left\{0, \frac{\pi}{4}, \frac{\pi}{2}, \dots\right\}.$$

1600 The transformed grasp poses are retained only if they remain collision-free and pass lightweight
 1601 grasp-stability checks.

1602 For a mug, symmetry-aware augmentation is more restrictive. Although the mug body may be
 1603 approximately cylindrical, the handle breaks rotational symmetry for functional grasping. Therefore,
 1604 a handle grasp is augmented by local perturbations on the handle, but it is not rotated around the mug

1605 body unless the transformed pose still lies on a functionally equivalent handle region. This prevents
1606 augmentation from turning a functional grasp-for-use candidate into a merely physical body grasp.

1607 For garment folding, local perturbation is more useful than global symmetry in many cases. Sleeve
1608 folding candidates are augmented by jittering the sleeve endpoint and mirrored place target within
1609 local confidence regions. Pants folding candidates are augmented by perturbing the centerline, side-
1610 fold target, and bottom-up fold target within bounded ranges. These perturbations create multiple
1611 folding trajectories with slightly different pickup and place points, which helps downstream policies
1612 learn robust behavior under keypoint noise and garment shape variation.

1613 For articulation, a cabinet handle candidate can be augmented by jittering the handle contact point
1614 and pull direction. If the cabinet has repeated equivalent handles, symmetry-aware or copy-based
1615 augmentation can transfer a validated handle grasp or articulation trajectory to another handle.
1616 However, each transferred candidate must be checked against the local collision context, since an
1617 adjacent wall, shelf, or drawer frame may make the copied trajectory infeasible.

1618 **Discussion.** Physics-aware augmentation increases annotation diversity while avoiding uncontrolled
1619 label noise. Local perturbation captures natural variation around a functional anchor, such as different
1620 grasp centers, fold points, pull directions, and base poses. Symmetry-aware augmentation exploits
1621 object-level or scene-level regularities, such as rotationally symmetric bottles or repeated cabinet
1622 handles. Both forms of augmentation are constrained by visual grounding, functional affordance, and
1623 physical feasibility checks. As a result, the expanded candidate banks contain diverse candidates that
1624 remain executable and functionally meaningful.

1625 **D Primitive-specific Action Annotation**

1626 **Parallel-jaw and dexterous grasp annotation.**

1627 **Articulation waypoint annotation.**

1628 **Insertion and hanging annotation.**

1629 **Garment and deformable-object trajectory annotation.**

1630 **Navigation target and approach-pose annotation.**

1631 **E Additional detail of Large-scale Robot Data Collection**

1632 **Overview.** Our large-scale data collection system consumes AnnotateAnything annotations as
1633 executable task interfaces, rather than relying on manually scripted object-specific demonstrations.

1634 **Atomic-skill library.** We implement an annotation-aligned atomic-skill library covering table-
1635 top manipulation, bimanual manipulation, whole-body control, humanoid control, dexterous-hand
1636 manipulation, and mobile manipulation.

1637 **Skill composition.** Each atomic skill maps compact annotation fields, such as grasp poses, target
1638 parts, waypoints, insertion directions, hanging anchors, and navigation goals, into controller or
1639 planner goals.

1640 **Long-horizon tasks.** Long-horizon tasks are constructed by composing multiple atomic skills
1641 according to task templates and object-state transitions.

1642 **Heterogeneous asynchronous environments.** Rollouts are generated in asynchronous parallel
1643 simulation environments, where each worker independently samples assets, tasks, robot embodiments,
1644 object poses, obstacles, and scene layouts.

1645 **Motion planning backend.** For each rollout, we use cuRobo-v2 for goal-set IK, motion planning,
1646 and obstacle-aware execution, since it supports GPU-native collision-aware motion generation for
1647 both standard manipulators and high-DoF embodiments [36].

1648 **Domain randomization.** We apply domain randomization over object pose, robot-object configura-
1649 tion, camera viewpoint, lighting, material, texture, distractor objects, and scene layout.

1650 **Candidate selection.** Because randomization changes the relative pose between robot and object,
1651 each worker retrieves multiple candidates from the annotation bank, solves goal-set IK, and executes
1652 the feasible candidate with the lowest planning cost.

1653 **Trajectory validation.** After execution, we validate task success, collision safety, contact consis-
1654 tency, articulation progress, and final-state stability when applicable.

1655 **Candidate filtering.** Annotation candidates that repeatedly fail under randomized simulation are
1656 removed from the candidate bank to improve future rollout quality.

1657 **Scope.** This data-collection system is used here to demonstrate that AnnotateAnything annotations
1658 are executable and scalable, while a full system-level analysis is left to our companion work.

1659 E.1 Keypoint Generation and Affordance Grounding

1660 **Overview.** AnnotateAnything naturally supports two related but different perception labels: key-
1661 points and affordances. Keypoints are produced by the visual annotation stage, while affordances
1662 are derived from physics-validated interaction annotations. The distinction is important: keypoints
1663 describe sparse semantic or functional anchors on an object or scene, whereas affordances describe
1664 where and how a robot can execute an interaction. Thus, in our framework, keypoint generation
1665 is mainly a visual grounding problem, while affordance grounding is an executable-interaction
1666 grounding problem.

1667 **Keypoint generation.** As described in Sec. 3.1.2, we generate keypoint annotations from fused 3D
1668 observations of the asset. Given a rendered RGB-D sequence, we first reconstruct a point cloud and
1669 sample diverse candidate points, e.g., using farthest point sampling. A VLM then selects points that
1670 are semantically or functionally important, such as handles, openings, corners, garment landmarks, or
1671 scene-level target locations. These selected keypoints provide sparse 3D anchors for downstream
1672 perception and action generation. They are not required to be directly executable by themselves;
1673 instead, they identify meaningful regions that later physics-based action annotation can refine and
1674 validate.

1675 **Affordance grounding from validated interactions.** Affordance labels are generated from the
1676 physics-based action annotation pipeline in Sec. 3.2. Given visual anchors, the action pipeline
1677 generates candidate interactions and validates them through geometry checks, motion planning,
1678 collision checking, contact reasoning, and physics simulation. We therefore use the validated
1679 candidate bank as affordance supervision. Compared with manually annotated affordance ground
1680 truth, which often marks regions that are visually or semantically plausible, our labels indicate regions
1681 and directions that are executable under physical constraints. For robot learning, this can be a stronger
1682 supervision signal: the model is trained to predict not only where an interaction appears possible, but
1683 where it has been verified to succeed.

1684 **Why physics-validated affordances are useful.** Manual affordance labels are usually subjective,
1685 sparse, and task-agnostic. For example, a human annotator may mark the handle of a mug or drawer
1686 as graspable, but the label may not specify which grasp pose, approach direction, contact mode, or
1687 motion trajectory is feasible for a robot. In contrast, our validated interaction annotations encode
1688 the full action context, including the target anchor, skill type, action direction, contact configuration,
1689 trajectory, and success metadata. They also preserve multiple valid solutions instead of collapsing the
1690 affordance to a single canonical point. This is especially important for manipulation, where many
1691 grasps, pulls, placements, or navigation poses can be valid for the same object.

1692 **Rigid objects.** For rigid objects, affordance labels are represented by an interaction anchor together
1693 with action parameters from validated grasping or dexterous-grasping candidates. For parallel-jaw
1694 grasping, the label may include the grasp center, approach direction, gripper orientation, and width.
1695 For dexterous hands, the label may additionally include deep palm pose, finger contact locations, and

1696 contact mode. These labels can be projected onto object point clouds or RGB-D images as affordance
1697 heatmaps, while the approach direction or grasp frame can be used as directional supervision.

1698 **Articulated objects.** For articulated objects, affordance grounding includes both the contact anchor
1699 and the motion trajectory. For example, opening a drawer requires identifying the handle anchor
1700 and predicting a pulling trajectory, while rotating a door or knob requires an anchor together with a
1701 rotation-aware motion path. Thus, the affordance label is not only “where to touch”, but also “how
1702 the interaction should move”. This allows the same object part to support different task-conditioned
1703 affordances, such as pulling, pushing, rotating, or holding.

1704 **Garments and deformable objects.** For garments and deformable objects, affordances are repre-
1705 sented by pick-and-place keypoints and their associated manipulation trajectories. Examples include
1706 bimanual pick points for lifting or spreading, pick-place pairs for folding, and hanging points for
1707 placing garments on a support. Unlike rigid-object affordances, garment affordances are often rela-
1708 tional: a pick point is meaningful only together with its paired place point, second-hand grasp point,
1709 or target folding line. We therefore represent garment affordance labels as structured point pairs or
1710 point sets rather than isolated single-point heatmaps.

1711 **Room-scale scenes and navigation affordances.** For room-scale scenes, affordances are grounded
1712 on occupancy maps, BEV maps, and scene layouts. The labels include navigable regions, object-
1713 centric approach poses, interaction-ready base poses, and feasible navigation targets. These room-
1714 level affordances are useful for mobile manipulation because the robot must first reach a suitable base
1715 pose before executing object-level manipulation. Thus, scene affordance grounding connects global
1716 navigation with local interaction.

1717 **Projection to perception supervision.** The generated labels can be converted into different super-
1718 vision formats depending on the downstream model. For image-based models, 3D anchors, contact
1719 regions, and trajectories are projected into rendered RGB-D views to form keypoint targets, affordance
1720 masks, heatmaps, or directional fields. For point-cloud models, the same labels are attached directly
1721 to 3D points or local surface patches. For navigation models, room-level affordances are rasterized
1722 into BEV maps with traversability, target-pose, and interaction-readiness labels. When multiple
1723 validated candidates exist for the same anchor or skill, we keep the multi-modal label distribution
1724 rather than forcing a single ground-truth target.

1725 **Positive and negative supervision.** Validated candidates provide positive affordance labels. Failed
1726 candidates can also be used as hard negatives when their failure is caused by physical infeasibility,
1727 such as collision, unstable contact, unreachable geometry, or invalid articulation motion. This
1728 provides richer supervision than manual affordance masks, since the model can learn the boundary
1729 between visually plausible but physically invalid regions and truly executable interaction regions.
1730 In practice, we treat negative labels conservatively, since some failures may depend on a particular
1731 embodiment, planner, or randomized scene configuration rather than the intrinsic affordance of the
1732 object.

1733 **Scope.** We use this construction mainly to demonstrate that AnnotateAnything annotations can
1734 be repurposed as robot-centric perception supervision. A full benchmark of affordance or keypoint
1735 detectors is outside the main scope of this paper. The key point is that our annotation format already
1736 contains the necessary supervision: visual annotations provide semantic 3D keypoints, and physics-
1737 validated action annotations provide executable affordance labels across rigid objects, articulated
1738 objects, garments, and room-scale scenes.

1739 E.2 Robot Reasoning VQA

1740 **Overview.** Our trajectory generation pipeline records execution-level metadata during simulation
1741 rollout, which can be reused to construct robot reasoning VQA pairs. Unlike static asset-level QA,
1742 these questions focus on what the robot actually executes under a sampled scene, robot pose, and
1743 object configuration.

1744 **Runtime-grounded supervision.** For each rollout, we record structured information such as the
1745 selected object, target part, anchor point, skill type, IK goal-set solution, approach pose, navigation

1746 target, articulation direction, garment manipulation point, and task success state. These labels are
1747 determined at runtime because domain randomization changes the relative pose between the robot
1748 and the object. Thus, even for the same asset, the executed candidate may vary across rollouts.

1749 **QA construction.** We convert the recorded trajectory metadata into question-answer pairs with
1750 VLM-assisted language generation. The VLM is used only to diversify and naturalize the question
1751 wording, while the answers are derived from simulator records and selected action candidates. This
1752 avoids treating the VLM prediction itself as ground truth.

1753 **Question categories.** The resulting QA pairs cover several robot-centric reasoning types, including
1754 object and part selection, approach direction, navigation progress, articulation motion, garment
1755 manipulation landmarks, and execution success. For example, questions may ask which object or
1756 part the robot is interacting with, from which side the robot approaches a target, where the robot
1757 is navigating, which direction an articulated part is moved, or which garment landmark is selected
1758 for manipulation. These questions require trajectory-level execution traces and cannot be reliably
1759 answered from static asset annotations alone.

1760 **Scope.** We use Robot Reasoning VQA as a lightweight downstream task to show that
1761 AnnotateAnything-generated rollouts naturally provide grounded multimodal reasoning supervi-
1762 sion. A full VQA benchmark is outside the main scope of this paper.

1763 E.3 3D VLM Instruction-tuning Data

1764 **Overview.** AnnotateAnything can be used to construct lightweight instruction-tuning data for 3D
1765 VLMs. This data comes from two sources: simulation-grounded visual labels and language-based
1766 scene annotations. The goal is not to introduce a new 3D VLM benchmark, but to show that the
1767 generated annotations can be naturally repurposed into multimodal instruction data.

1768 **Simulation-grounded visual supervision.** Because our assets and trajectories are generated in
1769 simulation, we can obtain dense ground-truth visual labels without manual annotation. These include
1770 3D bounding boxes, projected 2D bounding boxes, instance segmentation masks, object poses,
1771 part identities, and object-instance associations. Such labels can be converted into grounding-style
1772 instruction-response pairs, where the model is asked to localize objects, identify instances, refer to
1773 parts, or associate 2D observations with 3D scene elements. The answers are derived from simulator
1774 metadata rather than from VLM predictions.

1775 **Language and spatial-reasoning supervision.** The language annotations and cross-level compo-
1776 sition in Sec. 3.1.3 provide another source of instruction data. Since object-level annotations are
1777 linked to room-level layouts through instance identities and 6D poses, we can generate QA pairs
1778 about object relations, room structure, scene layout, and task context. Typical questions involve 3D
1779 spatial relationships, such as relative position, containment, support, proximity, and object-to-room
1780 associations. This makes the data suitable for training models to reason over both object-level
1781 semantics and scene-level spatial structure.

1782 **Instruction construction.** We convert structured annotations into instruction-response pairs using
1783 templates and VLM-assisted paraphrasing. Templates ensure that each answer is grounded in simula-
1784 tor or annotation metadata, while paraphrasing improves linguistic diversity. For example, grounding
1785 instructions may ask the model to locate an object or part, while spatial-reasoning instructions may
1786 ask about the relation between two objects or the functional area of a room. We avoid using the VLM
1787 as the source of ground truth; it is only used to diversify the language form.

1788 **Scope.** This downstream task mainly demonstrates that AnnotateAnything provides reusable super-
1789 vision for 3D multimodal learning. The same annotated assets can produce visual grounding labels
1790 from simulation and spatial-reasoning QA pairs from language and cross-level scene composition. A
1791 full-scale 3D VLM training and evaluation study is left outside the main scope of this paper.

1792 F Related Work

1793 F.1 Simulation-Based Automatic Data Collection

1794 Prior work on simulation-based data generation can be divided into two broad families. The first family
1795 is annotation- or demonstration-based. RLBench generates large-scale expert demonstrations through
1796 waypoint-guided motion planning across many tasks [59]. MimicGen amplifies a small number of
1797 human demonstrations by recomposing subtask segments across new scenes, object instances, and
1798 robot embodiments [60]. More recent systems such as RoboGen and GenSim use foundation models
1799 to propose tasks, synthesize scenes, and generate supervision or task code automatically [61, 62].
1800 These methods are attractive because they preserve a relatively interpretable supervision signal and
1801 often produce demonstrations that are closer to human decomposition of manipulation tasks.

1802 The second family is RL-based automatic data collection. Here, the simulator itself is the generator:
1803 environments such as Meta-World and Isaac Gym scale data through massive parallel interaction,
1804 while dexterous manipulation systems trained entirely in simulation demonstrate that difficult ma-
1805 nipulation can be learned without requiring human demonstrations [63–65]. The advantage is scale
1806 through optimization, but the cost is high engineering burden in reward design, curriculum design, and
1807 training stability, and the resulting supervision is often controller-centric rather than human-readable
1808 or reusable as annotation.

1809 A persistent gap across both families is that annotation-based systems usually depend on either simple
1810 tasks, a limited library of scripted decompositions, or scarce seed demonstrations, while RL-based
1811 systems can scale interaction but rarely produce portable semantic annotations that transfer across
1812 embodiments and task families. In other words, prior systems either generate *data for a fixed training*
1813 *setup* or *policies for a fixed task formulation*, rather than a reusable annotation layer that can support
1814 many downstream learners.

1815 **Connection to AnnotateAnything.** AnnotateAnything is motivated by exactly this gap: instead of
1816 choosing between narrow annotation reuse and reward-heavy RL optimization, it aims to generate
1817 reusable manipulation annotations that can bootstrap many downstream data-generation and policy-
1818 learning pipelines.

1819 F.2 Automatic Manipulation Annotation and Affordance Generation

1820 Automatic manipulation annotation has been studied extensively in grasping. Dex-Net 2.0 showed
1821 that robust grasp supervision can be synthesized at scale with analytic metrics in simulation [66].
1822 ACRONYM further expanded large-scale simulation-based grasp labels [67], while 6-DOF GraspNet
1823 and Contact-GraspNet demonstrated learned grasp proposal generation directly from point clouds and
1824 cluttered observations [68, 69]. For dexterous manipulation, DexGraspNet substantially increased
1825 the availability of robot-hand grasp labels over diverse objects [70]. Very recent work also begins to
1826 address bimanual dexterous grasp synthesis, but this remains significantly less mature than single-
1827 handed and parallel-gripper grasp annotation [71].

1828 Beyond grasping, affordance generation has been explored for articulated interaction, object-object
1829 relations, and task-specific manipulation. Where2Act and VAT-Mart learn where and how to interact
1830 with articulated objects from simulated interaction data [72, 73]. GPartNet argues that generalizable
1831 and actionable parts are a better transfer unit than object categories alone, providing part semantics and
1832 part poses that bridge perception and manipulation [74]. O2O-Afford studies object-object affordance
1833 without manual labels, covering behaviors such as placement and fitting [75]. OmniHang addresses
1834 hanging as a dedicated contact- and correspondence-centric problem [76], while assembly systems
1835 such as IndustReal focus on insertion and contact-rich assembly through task-specific learning and
1836 simulation design [77]. 3D AffordanceNet provides a benchmark for visual affordance understanding
1837 in 3D, but still within a fixed affordance ontology [78].

1838 Taken together, these works reveal a recurring tradeoff. Optimization-based annotation pipelines
1839 often provide physically grounded labels but are usually tied to a specific manipulation family, gripper
1840 model, contact formulation, or embodiment. Learned affordance models expand the label space
1841 and can capture richer semantics, but most are trained for a narrow set of task families, a narrow
1842 embodiment range, or a fixed affordance vocabulary. As a result, the field contains many strong

1843 *single-skill annotation engines*, but very few systems that can annotate many manipulation skills
1844 across many embodiments in a unified way.

1845 **Connection to AnnotateAnything.** AnnotateAnything is best positioned as a general annotation
1846 substrate above task-specialized affordance systems: it seeks to unify grasp-, articulation-, placement-,
1847 insertion-, and hanging-style supervision under one scalable pipeline rather than introducing one
1848 more narrow annotation engine.

1849 **F.3 Large-Scale 3D Assets and Generative Asset Creation**

1850 Large-scale 3D asset repositories underpin almost all recent simulation-based manipulation work.
1851 ShapeNet established the modern large-scale CAD repository paradigm [79], PartNet added fine-
1852 grained and hierarchical part annotations [80], and PartNet-Mobility extended the ecosystem with
1853 articulated joints and motion metadata through the SAPIEN dataset stack [81, 82]. Objaverse
1854 dramatically increased asset scale and diversity by aggregating hundreds of thousands of annotated
1855 3D objects from the web [83]. For manipulation-specific repositories, ACRONYM and DexGraspNet
1856 provide grasp annotations, while GAPartNet adds action-relevant part semantics and poses [67, 70,
1857 74].

1858 A parallel thread focuses on generative asset creation. DreamFusion and Shap-E demonstrate that
1859 text-conditioned 3D generation can produce novel 3D assets without manual modeling [84, 85].
1860 These methods are promising for simulation scale-up, but they usually do not provide the robot-ready
1861 metadata that manipulation systems need: articulated joints, functional parts, interaction affordances,
1862 contact semantics, or embodiment-aware trajectories.

1863 The central gap is therefore not simply the amount of 3D content. The field already has abundant
1864 static geometry and increasingly capable generative models. What remains scarce is a pipeline that
1865 can transform raw or generated assets into *interaction-ready annotations* suitable for manipulation
1866 learning, benchmark construction, and skill composition.

1867 **Connection to AnnotateAnything.** AnnotateAnything can be framed as the missing conversion
1868 layer from assets to actionable supervision: it complements repositories and text-to-3D generators by
1869 adding the manipulation semantics that raw geometry alone does not provide.

1870 **F.4 Simulation Environments and Manipulation Benchmarks**

1871 Robot learning has benefited from a rich simulator ecosystem. MuJoCo remains a dominant engine
1872 for fast rigid-body control and benchmarked continuous-control research [86]. PyBullet offers an
1873 accessible open-source simulation stack with broad adoption in robotics and reinforcement learning
1874 practice [87]. Isaac Gym introduced end-to-end GPU-resident simulation for massive parallel robot
1875 learning [64], and Isaac Sim extends the Isaac stack toward richer industrial and photorealistic
1876 simulation workflows [88]. SAPIEN focuses on interactive articulated-object simulation and has
1877 become especially influential for part-aware and household manipulation research [82]. RoboSuite
1878 provides a modular MuJoCo-based research framework for manipulation [89].

1879 On top of these platforms, benchmark suites define the evaluation targets. Meta-World emphasizes
1880 multitask and meta-RL over tabletop manipulation skills [63]. RLBench provides a large task library
1881 with automatically generated demonstrations [59]. ManiSkill and ManiSkill2 move toward more
1882 diverse objects, richer demonstrations, and a broader range of manipulation modes and embodiments
1883 [90, 91]. BEHAVIOR broadens the focus to everyday household activities and long-horizon embodied
1884 tasks [92].

1885 These environments and benchmarks are foundational, but they expose a systems-level fragmentation.
1886 Some emphasize control fidelity, others articulated assets, others broad task coverage, and others long-
1887 horizon embodied realism. What they usually do *not* provide is a unified mechanism for generating
1888 transferable task annotations across assets, embodiments, and manipulation families. In practice, the
1889 benchmark and simulator literature still relies heavily on manual task authoring, limited affordance
1890 ontologies, or benchmark-specific annotation choices.

1891 **Connection to AnnotateAnything.** AnnotateAnything can be presented as complementary infras-
1892 tructure for this ecosystem: rather than replacing simulators or benchmarks, it aims to make them
1893 easier to populate with interaction-ready annotations at larger scale and with greater task diversity.

1894 **F.5 Skill-Centric Robot Learning and Atomic Manipulation Skills**

1895 A separate but highly relevant line of work studies how robots can compose reusable skills. The
1896 options framework is the classic formal foundation for temporally extended actions and hierarchical
1897 decision making [93]. In robotics, recent systems operationalize this idea with explicit or implicit
1898 skill libraries: SayCan selects among a predefined set of robotic skills by combining language-model
1899 usefulness with affordance-grounded feasibility [94], while Code as Policies composes control
1900 APIs and primitive calls through language-model-generated programs [95]. Composable Part-Based
1901 Manipulation further pushes toward reusable, object-part-centered manipulation abstractions [96].

1902 What all of these systems have in common is that they assume some skill inventory, primitive API, or
1903 part-level grounding already exists. In long-horizon robot learning, the bottleneck often shifts from
1904 policy optimization to *grounding*: how to obtain the contact points, manipulable regions, trajectories,
1905 constraints, or object-part correspondences that instantiate a skill on a new object. Skill composition
1906 is therefore only as scalable as the annotation layer beneath it.

1907 **Connection to AnnotateAnything.** AnnotateAnything directly supports skill-centric learning by
1908 making the atomic grounding of skills more scalable: if reusable annotations can be generated
1909 automatically, then skill libraries and hierarchical planners can be expanded far beyond hand-authored
1910 primitive sets.

1911 **NeurIPS Paper Checklist**

1912 The checklist is designed to encourage best practices for responsible machine learning research,
1913 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
1914 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
1915 follow the references and follow the (optional) supplemental material. The checklist does NOT count
1916 towards the page limit.

1917 Please read the checklist guidelines carefully for information on how to answer these questions. For
1918 each question in the checklist:

- 1919 • You should answer [Yes], [No], or [N/A].
- 1920 • [N/A] means either that the question is Not Applicable for that particular paper or the
1921 relevant information is Not Available.
- 1922 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

1923 **The checklist answers are an integral part of your paper submission.** They are visible to the
1924 reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it
1925 (after eventual revisions) with the final version of your paper, and its final version will be published
1926 with the paper.

1927 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
1928 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a
1929 proper justification is given (e.g., error bars are not reported because it would be too computationally
1930 expensive” or “we were unable to find the license for the dataset we used”). In general, answering
1931 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we
1932 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
1933 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
1934 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
1935 please point to the section(s) where related material for the question can be found.

1936 IMPORTANT, please:

- 1937 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”.**
- 1938 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 1939 • **Do not modify the questions and only use the provided macros for your answers.**

1940 **1. Claims**

1941 Question: Do the main claims made in the abstract and introduction accurately reflect the
1942 paper’s contributions and scope?

1943 Answer: **[TODO]**

1944 Justification: **[TODO]**

1945 Guidelines:

- 1946 • The answer [N/A] means that the abstract and introduction do not include the claims
1947 made in the paper.
- 1948 • The abstract and/or introduction should clearly state the claims made, including the
1949 contributions made in the paper and important assumptions and limitations. A [No] or
1950 [N/A] answer to this question will not be perceived well by the reviewers.
- 1951 • The claims made should match theoretical and experimental results, and reflect how
1952 much the results can be expected to generalize to other settings.
- 1953 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1954 are not attained by the paper.

1955 **2. Limitations**

1956 Question: Does the paper discuss the limitations of the work performed by the authors?

1957 Answer: **[TODO]**

1958 Justification: **[TODO]**

1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [TODO]

Justification: [TODO]

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

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

Answer: [TODO]

Justification: [TODO]

Guidelines:

- The answer [N/A] means that the paper does not include experiments.

2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063

- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [TODO]

Justification: [TODO]

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

- 2064 • Providing as much information as possible in supplemental material (appended to the
2065 paper) is recommended, but including URLs to data and code is permitted.

2066 6. Experimental setting/details

2067 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
2068 rameters, how they were chosen, type of optimizer) necessary to understand the results?

2069 Answer: **[TODO]**

2070 Justification: **[TODO]**

2071 Guidelines:

- 2072 • The answer [N/A] means that the paper does not include experiments.
- 2073 • The experimental setting should be presented in the core of the paper to a level of detail
2074 that is necessary to appreciate the results and make sense of them.
- 2075 • The full details can be provided either with the code, in appendix, or as supplemental
2076 material.

2077 7. Experiment statistical significance

2078 Question: Does the paper report error bars suitably and correctly defined or other appropriate
2079 information about the statistical significance of the experiments?

2080 Answer: **[TODO]**

2081 Justification: **[TODO]**

2082 Guidelines:

- 2083 • The answer [N/A] means that the paper does not include experiments.
- 2084 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
2085 intervals, or statistical significance tests, at least for the experiments that support the
2086 main claims of the paper.
- 2087 • The factors of variability that the error bars are capturing should be clearly stated (for
2088 example, train/test split, initialization, random drawing of some parameter, or overall
2089 run with given experimental conditions).
- 2090 • The method for calculating the error bars should be explained (closed form formula,
2091 call to a library function, bootstrap, etc.)
- 2092 • The assumptions made should be given (e.g., Normally distributed errors).
- 2093 • It should be clear whether the error bar is the standard deviation or the standard error
2094 of the mean.
- 2095 • It is OK to report 1-sigma error bars, but one should state it. The authors should
2096 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
2097 of Normality of errors is not verified.
- 2098 • For asymmetric distributions, the authors should be careful not to show in tables or
2099 figures symmetric error bars that would yield results that are out of range (e.g., negative
2100 error rates).
- 2101 • If error bars are reported in tables or plots, the authors should explain in the text how
2102 they were calculated and reference the corresponding figures or tables in the text.

2103 8. Experiments compute resources

2104 Question: For each experiment, does the paper provide sufficient information on the com-
2105 puter resources (type of compute workers, memory, time of execution) needed to reproduce
2106 the experiments?

2107 Answer: **[TODO]**

2108 Justification: **[TODO]**

2109 Guidelines:

- 2110 • The answer [N/A] means that the paper does not include experiments.
- 2111 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
2112 or cloud provider, including relevant memory and storage.
- 2113 • The paper should provide the amount of compute required for each of the individual
2114 experimental runs as well as estimate the total compute.

- 2115 • The paper should disclose whether the full research project required more compute
2116 than the experiments reported in the paper (e.g., preliminary or failed experiments that
2117 didn't make it into the paper).

2118 9. Code of ethics

2119 Question: Does the research conducted in the paper conform, in every respect, with the
2120 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

2121 Answer: **[TODO]**

2122 Justification: **[TODO]**

2123 Guidelines:

- 2124 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
2125 Ethics.
- 2126 • If the authors answer [No], they should explain the special circumstances that require a
2127 deviation from the Code of Ethics.
- 2128 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
2129 eration due to laws or regulations in their jurisdiction).

2130 10. Broader impacts

2131 Question: Does the paper discuss both potential positive societal impacts and negative
2132 societal impacts of the work performed?

2133 Answer: **[TODO]**

2134 Justification: **[TODO]**

2135 Guidelines:

- 2136 • The answer [N/A] means that there is no societal impact of the work performed.
- 2137 • If the authors answer [N/A] or [No], they should explain why their work has no societal
2138 impact or why the paper does not address societal impact.
- 2139 • Examples of negative societal impacts include potential malicious or unintended uses
2140 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
2141 (e.g., deployment of technologies that could make decisions that unfairly impact specific
2142 groups), privacy considerations, and security considerations.
- 2143 • The conference expects that many papers will be foundational research and not tied
2144 to particular applications, let alone deployments. However, if there is a direct path to
2145 any negative applications, the authors should point it out. For example, it is legitimate
2146 to point out that an improvement in the quality of generative models could be used to
2147 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
2148 that a generic algorithm for optimizing neural networks could enable people to train
2149 models that generate Deepfakes faster.
- 2150 • The authors should consider possible harms that could arise when the technology is
2151 being used as intended and functioning correctly, harms that could arise when the
2152 technology is being used as intended but gives incorrect results, and harms following
2153 from (intentional or unintentional) misuse of the technology.
- 2154 • If there are negative societal impacts, the authors could also discuss possible mitigation
2155 strategies (e.g., gated release of models, providing defenses in addition to attacks,
2156 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
2157 feedback over time, improving the efficiency and accessibility of ML).

2158 11. Safeguards

2159 Question: Does the paper describe safeguards that have been put in place for responsible
2160 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
2161 image generators, or scraped datasets)?

2162 Answer: **[TODO]**

2163 Justification: **[TODO]**

2164 Guidelines:

- 2165 • The answer [N/A] means that the paper poses no such risks.

- 2166
- 2167
- 2168
- 2169
- 2170
- 2171
- 2172
- 2173
- 2174
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

2175

2176

2177

2178

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

2179

Answer: **[TODO]**

2180

Justification: **[TODO]**

2181

Guidelines:

- 2182
- 2183
- 2184
- 2185
- 2186
- 2187
- 2188
- 2189
- 2190
- 2191
- 2192
- 2193
- 2194
- 2195
- 2196
- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

2197

2198

2199

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

2200

Answer: **[TODO]**

2201

Justification: **[TODO]**

2202

Guidelines:

- 2203
- 2204
- 2205
- 2206
- 2207
- 2208
- 2209
- 2210
- The answer [N/A] means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

2211

2212

2213

2214

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

2215

Answer: **[TODO]**

2216

Justification: **[TODO]**

2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

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

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

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

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.

Table 7: Example instantiations of the unified action annotation schema for a cabinet or drawer. The handle anchor carries the functional affordance of opening, while scene-level anchors support navigation and approach-pose annotation.

| Symbol | Cabinet example | Interpretation |
|---|------------------------|--|
| \mathcal{A} | Cabinet scene | The input asset or scene, including cabinet doors, drawers, handles, panels, and surrounding free space. |
| h_{handle} | Cabinet handle | A part-level anchor from $\mathcal{P}(\mathcal{A})$ and a functional affordance anchor. It can support both grasping and articulated manipulation. |
| $\phi(h_{\text{handle}})$ | Grasp-for-opening | The functional affordance of the handle. It indicates that grasps on this anchor should support opening or closing the door or drawer. |
| s_{grasp} | Parallel-jaw grasp | A compatible skill in $\mathcal{S}(h_{\text{handle}})$. This skill instantiates grasp poses for holding the cabinet or drawer handle. |
| $\mathcal{B}_{h_{\text{handle}}, s_{\text{grasp}}}$ | Handle grasp bank | A set of feasible grasp candidates on the handle, including different contact points, gripper poses, approach directions, and gripper widths. These candidates are conditioned on the handle’s role as an opening affordance. |
| $a_{h,s}^{(i)}$ | One handle grasp | One concrete grasp candidate. Here, $x^{(i)}$ is a contact point or patch on the handle, $\theta^{(i)}$ stores the 6D grasp pose and gripper width, $\tau^{(i)}$ stores the approach trajectory, $v^{(i)}$ records IK, collision, and functional feasibility, and $d^{(i)}$ records approach direction and grasp family. |
| s_{art} | Articulation | Another compatible skill in $\mathcal{S}(h_{\text{handle}})$. The same handle anchor can instantiate opening or closing actions for a door or drawer. |
| $\mathcal{B}_{h_{\text{handle}}, s_{\text{art}}}$ | Articulation bank | A set of feasible articulation candidates, including different pull directions, waypoint sequences, contact-maintenance strategies, and motion ranges. |
| $a_{h,s}^{(i)}$ | One opening trajectory | One concrete articulation candidate. Here, $x^{(i)}$ is the handle interaction point, $\theta^{(i)}$ stores the prismatic or revolute motion direction, $\tau^{(i)}$ contains pull or rotate waypoints, $v^{(i)}$ records collision status, contact maintenance, actuation success, and functional success, and $d^{(i)}$ records pull direction, waypoint spacing, and trajectory family. |
| h_{free} | Free-space region | A scene-level anchor from $\mathcal{R}^{\text{scene}}(\mathcal{A})$, corresponding to reachable space in front of the cabinet. |
| $\phi(h_{\text{free}})$ | Navigation approach | The functional affordance of the free-space region. It indicates that this region supports approaching the cabinet for subsequent manipulation. |
| s_{nav} | Navigation | A compatible skill in $\mathcal{S}(h_{\text{free}})$. This skill generates interaction-ready base poses and approach directions. |
| $\mathcal{B}_{h_{\text{free}}, s_{\text{nav}}}$ | Navigation bank | A set of feasible base-pose candidates from which the robot can approach and manipulate the cabinet, with different distances, orientations, viewpoints, and approach sides. |
| $a_{h,s}^{(i)}$ | One base pose | One concrete navigation candidate. Here, $x^{(i)}$ is a target base pose, $\theta^{(i)}$ stores base orientation and approach direction, $\tau^{(i)}$ stores an optional collision-free path, $v^{(i)}$ records path feasibility, manipulation reachability, and functional approach validity, and $d^{(i)}$ records viewpoint, clearance, and distance-to-object. |

Table 8: Examples of candidate region construction for different skills. Candidate regions are obtained from language annotations, visual grounding, and asset geometry before concrete target anchors are sampled.

| Skill | Annotation source | Candidate region | Example interpretation |
|----------------------|-----------------------------------|---|--|
| Functional grasp | Part mask + language | Mug handle, tool handle | The language annotation indicates that the handle is used for holding, pouring, or operating the object. Targets sampled from this region are functional grasp candidates. |
| Physical grasp | Part mask + geometry | Mug body, box side, object surface | The region is physically graspable but not necessarily functionally preferred. For example, grasping a mug body may be stable for pickup but less suitable for pouring. |
| Dexterous grasp | Part mask + affordance region | Handle, body surface, functional contact area | The region supports multi-finger contacts. Candidates differ in fingertip locations, palm poses, and hand configurations. |
| Garment manipulation | Semantic keypoints | Cloth corners, edges, sleeve endpoints | Keypoints directly seed grasp anchors or bimanual keypoint pairs for lifting, folding, stretching, or placing. |
| Articulation | Part mask + affordance region | Cabinet handle, drawer handle, door edge | The handle or movable part is selected as the region from which pulling, pushing, or rotating anchors are sampled. |
| Insertion | Affordance region + geometry | Opening, socket, hole, container mouth | The opening region provides entry points, alignment frames, and insertion directions. |
| Hanging | Affordance region + edge geometry | Hook, support edge, handle loop | Support regions are sampled to generate hanging contacts and gravity-stable poses. |
| Navigation | Occupancy / BEV map | Reachable free-space region | Traversable cells around a target object are sampled to generate interaction-ready base poses and approach directions. |

Table 9: Examples of sampling and filtering rules for candidate target generation. These filters convert grounded candidate regions into feasible target anchors for candidate banks.

| Skill | Sampling method | Filtering criteria | Example rejected targets |
|----------------------|---|---|---|
| Functional grasp | FPS on functional part mask | Part size, gripper width, surface normal, collision margin, reachability, functional affordance | A tiny handle decoration that is too small to grasp; a mug-body grasp when the task requires pouring from the handle. |
| Physical grasp | FPS on graspable surfaces | Contact area, antipodal geometry, local curvature, gripper clearance, IK feasibility | A highly curved or sharp surface patch that cannot support stable contact. |
| Dexterous grasp | FPS or contact-set sampling on part masks | Hand joint limits, fingertip support, self-collision, object collision, contact stability | A contact set that requires impossible finger extension or causes palm-object penetration. |
| Garment manipulation | Keypoint sampling and perturbation | Visibility, bimanual reachability, corner confidence, pair distance, deformation feasibility | Two cloth keypoints that are too far apart for bimanual grasping or are occluded. |
| Articulation | FPS on handle or movable-part mask | Motion-axis consistency, contact maintenance, collision, reachability, actuation direction | A handle point that can be grasped but cannot be pulled along the door or drawer motion direction. |
| Insertion | Sampling on opening region | Opening size, entry clearance, alignment direction, collision-free approach, insertion depth | An opening that is detected visually but too narrow for the inserted object. |
| Hanging | Sampling on hooks or support edges | Edge orientation, gravity stability, clearance, collision, contact support | A support edge that faces downward or lacks clearance for the hanging object. |
| Navigation | Grid sampling or 2D FPS on BEV free space | Traversability, clearance, path connectivity, visibility, manipulation reachability | A base pose that is collision-free but too far to reach the cabinet handle. |

Table 10: Trajectory generation families grouped by the source of constraints. Each generator converts a sampled target anchor into an initial trajectory or waypoint sequence.

| Constraint family | Representative skills | Main inputs | Generated trajectory |
|-------------------------------|--|--|---|
| Object-property-conditioned | Grasp, dexterous grasp, garment fling, garment fold | Object geometry, part masks, keypoints, garment length, grasp frame | Pre-grasp, contact, closure, retreat, bimanual garment motion, or keypoint-driven folding trajectory. |
| Object-trajectory-conditioned | Open, close, rotate, push articulated parts | Articulation axis, joint state, moving-part transform, handle pose | End-effector waypoints that follow the moving part while preserving grasp or contact constraints. |
| Object-vector-conditioned | Place, hang, pour, insert | Surface normal, insertion axis, hanging direction, pouring direction, active/passive vectors | Approach and execution waypoints aligned with a semantic or geometric vector. |
| Scene-trajectory-conditioned | Navigation, approach-pose generation, retrieval-style motion | Occupancy map, BEV map, traversable free space, target object region | Base-pose sequence, collision-free approach path, and interaction-ready base trajectory. |

Table 11: Garment trajectory templates generated from semantic keypoints. These trajectories are template-level seeds and are later optimized and validated.

| Garment skill | Keypoints | Geometric rule | Generated trajectory |
|--------------------|--|--|---|
| Sleeve fold | Sleeve endpoint, shoulder, bottom corner | Use shoulder–bottom line as fold line; mirror sleeve endpoint inward. | Pick sleeve endpoint, lift, move to mirrored target, place inside garment body. |
| Bottom-up fold | Bottom corners, shoulder points | Move bottom keypoints toward the shoulder line. | Bimanual pick of bottom corners, lift, fold upward toward shoulder region, place. |
| Pants side fold | Left/right leg or side keypoints, central axis | Fold the two side or leg regions toward the centerline. | Pick side or leg keypoints, lift, move inward, place to overlap legs. |
| Pants compact fold | Cuff/bottom keypoints, waistband region | Fold bottom or cuff region upward toward waistband; optionally repeat. | Pick bottom/cuff keypoints, fold upward, then perform a second compacting fold. |
| Fling | Two garment keypoints, garment length | Scale lift height and stretch distance by keypoint distance or garment length. | Bimanual lift, stretch, optional shake, and release or place trajectory. |

Table 12: Trajectory templates for articulated-object actions. These templates differ mainly in the joint schedule, contact mode, and whether the end-effector maintains a grasp or pushes on the part.

| Action type | Subcase | Trajectory rule |
|---------------------|------------------------------|---|
| Open with handle | Handle open | Grasp the handle, preserve T_{rel} , and increase the revolute or prismatic joint state from closed to open. |
| Partially open pull | Partial pull | Grasp the handle and follow the same relative-transform rule, but stop at a partial target joint state $q_{partial}$. |
| Close with handle | Handle close | Grasp the handle and reverse the joint schedule, moving from the current joint state toward the closed state. |
| Push close | Collision close / push close | Contact the door, drawer, or lid surface without necessarily forming a handle grasp; push along the closing direction while maintaining collision-free contact. |
| Open lid | Lid open | Follow the lid’s revolute joint trajectory while preserving the end-effector–lid relation or maintaining a pushing contact. |
| Close lid | Lid close | Reverse the lid joint trajectory and push or guide the lid toward the closed state. |
| Rotate button | Button rotate | Generate rotational waypoints around the button axis, with the end-effector orientation following the required rotation. |
| Rotate lid | Lid rotate | Rotate the lid around its axis, using angular waypoints and contact-maintenance constraints. |
| Push button | Button push | Generate a short linear trajectory along the button normal or push direction, with displacement bounded by the button travel range. |

Table 13: Object-vector-conditioned trajectory templates. These actions are generated by aligning the end-effector or object with a semantic or geometric vector.

| Skill | Vector source | Alignment rule | Generated trajectory |
|--------|---|---|---|
| Place | Surface normal and pose range | Align object placement pose with support normal. | Approach placement pose, lower along normal, release, and retreat. |
| Hang | Support edge, hook direction, gravity | Align hanging contact with support edge and gravity-stable direction. | Approach hook or edge, establish contact, lower to stable hanging pose, release. |
| Pour | Active pour vector and receiving direction | Rotate object so active pour vector points toward target container or receiving region. | Lift, rotate around grasp frame, pour for a specified angular range, return or retreat. |
| Insert | Active insertion axis and passive socket axis | Align active axis with passive axis and translate along insertion direction. | Approach socket, align, move along insertion axis, optionally retreat or release. |

Table 14: Trajectory optimization objectives for different skill families. Optimization refines the template trajectory produced by the previous stage into an asset-specific and embodiment-aware candidate.

| Skill family | Optimized variables | Main objectives | Examples |
|---------------------------|---|---|--|
| Dexterous/contact skills | Hand pose, finger joints, contact assignment, approach pose | Contact alignment, functional-region consistency, wrench support, joint limits, self-collision, hand-object collision | Dexterous grasp on a mug handle, whole-hand grasp, bimanual grasp. |
| Adaptive object skills | Waypoint positions, placement targets, timing, approach offsets | Geometry adaptation, keypoint consistency, length/scale adaptation, bimanual reachability, smoothness | Garment sleeve fold, bottom-up fold, pants fold, fling. |
| Articulated-object skills | End-effector waypoints, contact pose, joint schedule, pull/push direction | Joint-motion tracking, contact maintenance, collision avoidance, reachability, smoothness | Open/close cabinet, pull drawer, rotate button, push button. |
| Vector-conditioned skills | Approach pose, alignment pose, execution direction, release pose | Vector alignment, insertion depth, support stability, placement normal consistency, collision-free approach | Insert, hang, pour, place. |
| Scene-level skills | Base path, velocity commands, base orientation, approach pose | Path following, clearance, traversability, local dynamics, turning radius, manipulation reachability | Navigation to cabinet, approach-pose generation, room-scale interaction. |

Table 15: Examples of articulation trajectory refinement. The optimizer adapts template waypoints to contact, collision, reachability, and joint-motion constraints.

| Action | Optimization focus | Typical adjustment |
|---------------------|---|---|
| Open with handle | Maintain grasp while following revolute or prismatic motion | Adjust pull direction, waypoint spacing, and wrist orientation to preserve contact and avoid collision. |
| Partially open pull | Stop at a partial target state | Shorten the joint schedule while preserving a feasible pull trajectory. |
| Close with handle | Reverse opening trajectory | Adapt closing waypoints to avoid collision with the cabinet frame and maintain grasp. |
| Push close | Maintain pushing contact without grasp | Align push direction with closing motion and regulate contact force and clearance. |
| Open/close lid | Follow lid rotation | Adjust end-effector orientation and approach side to avoid collision with the lid and surrounding object. |
| Rotate button/lid | Follow angular motion | Refine angular waypoints and wrist orientation around the rotation axis. |
| Push button | Linear push along button normal | Bound displacement by travel range and avoid side contact. |

Table 16: Examples of vector-conditioned trajectory optimization. These skills are refined by improving vector alignment, clearance, depth, and stability.

| Skill | Optimized variables | Objective | Example adjustment |
|--------|---|--|---|
| Insert | Alignment pose, insertion depth, approach offset | Axis alignment, entry clearance, collision-free insertion | Rotate the object to align with the socket axis and reduce insertion depth if collision occurs. |
| Hang | Approach pose, support contact, release pose | Support-edge alignment, gravity stability, clearance | Shift the placement pose so the object lowers onto the hook without colliding. |
| Pour | Wrist rotation, pouring angle, target direction | Pour-vector alignment, wrist feasibility, collision avoidance | Adjust pouring angle so the liquid direction points toward the receiving container. |
| Place | Placement pose, release height, retreat direction | Surface-normal alignment, stable placement, collision-free retreat | Lower along the support normal and retreat without sweeping nearby objects. |

Table 17: Embodiment-aware navigation optimization. The A* path from trajectory generation is refined by DWB-style local rollout using the robot’s motion model.

| Embodiment | Control model | Main constraints | Optimization effect |
|----------------------------------|--------------------------------------|--|---|
| Differential drive | (v, ω) | Velocity bounds, angular velocity bounds, clearance, path tracking | Smoothly follows the A* path while respecting coupled forward and rotational motion. |
| Ackermann drive | (v, δ) | Steering limit, minimum turning radius, acceleration bound, clearance | Removes sharp A* turns that are infeasible for car-like platforms. |
| Holonomic base | (v_x, v_y, ω) | Velocity bounds, lateral motion limits, clearance, manipulation reachability | Uses lateral motion to approach interaction targets from better viewpoints or reachability zones. |
| Quadruped / humanoid abstraction | Holonomic or quasi-holonomic command | Traversability, body clearance, foothold or stance feasibility if available | Produces approach poses that remain compatible with whole-body navigation and manipulation. |

Table 18: Metadata stored after trajectory optimization. These fields support candidate ranking, debugging, validation, and downstream sampling.

| Metadata field | Meaning | Example |
|------------------------|--|---|
| Optimized parameters | Final skill-specific parameters $\theta_*^{(i)}$ | Dexterous hand pose and joints, insertion alignment pose, hanging pose, base orientation. |
| Optimized trajectory | Final waypoint sequence $\tau_*^{(i)}$ | Approach-contact-retreat path, articulated opening waypoints, DWB-refined navigation path. |
| Objective breakdown | Individual optimization terms | Collision cost, contact cost, reachability cost, smoothness cost, functional consistency cost. |
| Embodiment feasibility | Whether the optimized candidate satisfies robot-specific constraints | Joint-limit status, turning-radius feasibility, pre-grasp reachability, local base dynamics feasibility. |
| Functional consistency | Whether the optimized candidate preserves $\phi(h)$ | Mug-handle grasp remains grasp-for-use; cabinet-handle grasp remains grasp-for-opening. |
| Diversity descriptor | How this candidate differs from others | Different approach direction, contact set, fold target, base-pose side, navigation route, or trajectory family. |

Table 19: Validation embodiments used for different skill families. Floating validation checks object-centric physical feasibility, while non-floating validation checks embodiment-dependent feasibility.

| Validation setting | Representative skills | What is checked | Motivation |
|-------------------------|--|--|---|
| Floating gripper | Parallel-jaw grasp, insertion, hanging, placement, handle articulation | Contact, collision, penetration, stable grasp, task progress, release stability | Validates local object interaction without over-constraining the annotation to a particular arm pose. |
| Floating dexterous hand | Dexterous grasp, dexterous handle grasp, dexterous articulation | Finger-object contact, hand-object collision, hand joint validity, object stability, contact maintenance | Tests whether a dexterous contact configuration actually constrains the object before downstream execution. |
| Full manipulator | Garment push, wash, retrieval, entangled deformable-object skills | IK, joint limits, singularities, arm collision, self-collision, garment-arm entanglement, deformable outcome | Needed when feasibility depends on the arm and body, not only local hand-object contact. |
| Base-aware validation | Navigation, approach-pose annotation, mobile manipulation setup | Traversability, local dynamics, clearance, turning radius, path connectivity, manipulation reachability | Ensures that the sampled base pose and path are feasible for the target robot base. |

Table 20: Robustness tests used during physics validation. Gravity randomization is mainly applied to parallel-jaw and dexterous grasp candidates.

| Robustness test | Applied skills | Purpose |
|-----------------------------|--|---|
| Random gravity directions | Parallel-jaw grasp, dexterous grasp | Reject grasps that only succeed under one gravity direction but fail when the load direction changes. |
| Gravity magnitude scaling | Parallel-jaw grasp, dexterous grasp | Use $1.5g_0$ to create a stricter stability test for weak or marginal contacts. |
| Object pose jitter | Grasp, dexterous grasp, insertion, hanging | Reject candidates that require extremely precise initialization or have no robustness margin. |
| Hand or gripper pose jitter | Grasp, dexterous grasp, articulation | Test whether small execution errors cause collision, contact loss, or task failure. |
| External perturbation force | Grasp, dexterous grasp, hanging, placement | Check whether the object remains stable under small disturbances after contact or release. |
| Trajectory perturbation | Garment, articulation, vector-conditioned skills | Detect candidates that are overly sensitive to waypoint timing, contact pose, or approach direction. |

Table 21: Skill-specific validation criteria. A candidate must satisfy both generic physics checks and task-specific success conditions.

| Skill family | Validation criteria | Typical success signal |
|----------------------|---|--|
| Parallel-jaw grasp | Approach collision, closure contact, object stability, gravity robustness, disturbance robustness | Object remains held under randomized gravity and perturbation. |
| Dexterous grasp | Joint validity, self-collision, hand-object penetration, contact coverage, object stability, gravity robustness | Object is actually constrained by the hand and does not slip or fall. |
| Articulation | Contact maintenance, joint progress, collision, penetration, handle blockage, grasp or push validity | Door, drawer, lid, or button reaches the intended joint state without contact loss or jamming. |
| Insertion | Axis alignment, entry clearance, collision-free approach, insertion depth, post-insertion stability | Object reaches the target opening or socket without collision and remains aligned. |
| Hanging | Support contact, gravity stability, release outcome, clearance around hook or edge | Object remains hanging after release and does not slide or fall. |
| Garment / deformable | Full-arm IK, singularity avoidance, arm-cloth collision, entanglement, geometric outcome | Cloth reaches the intended fold, push, wash, or retrieval state without wrapping around the arm. |
| Navigation | Traversability, clearance, local dynamics, path connectivity, final-pose reachability | Robot reaches an interaction-ready base pose from which manipulation is feasible. |

Table 22: Common failure cases detected by physics validation. These failures motivate the use of simulation-based filtering instead of relying only on geometric candidate generation.

| Skill | Failure case | Cause | How validation detects it |
|------------------------------------|--|--|--|
| Parallel-jaw grasp | Object drops under different gravity directions | The gripper has too little contact area, contacts only an edge, or does not form a stable antipodal grasp. | The grasp may pass nominal lifting but fails under one of the eight randomized gravity directions or disturbance rollouts. |
| Parallel-jaw grasp | Immediate gripper-object collision | The pre-grasp or approach pose intersects the object or nearby geometry. | Collision and penetration statistics exceeded the tolerance before closure. |
| Dexterous grasp | Looks like a grasp but does not actually hold the object | The hand pose surrounds the object visually, but fingertips do not establish effective contacts or contact forces. | The object slips or falls during gravity robustness tests; contact coverage or stability score is too low. |
| Dexterous grasp | Finger penetration or self-collision | The optimized joint configuration places fingers inside the object or causes hand self-intersection. | Hand-object penetration, self-collision, or joint-limit violation is detected during rollout. |
| Dexterous articulation | Hand gets stuck in the handle or part geometry | Fingers enter a narrow handle, cavity, or gap and cannot move with the articulated part. | Contact maintenance fails, joint motion stalls, or penetration grows during opening or closing. |
| Dexterous articulation | Dexterous hand joint deformation or damage in simulation | The moving object part pushes against fingers or joints in an infeasible configuration. | Joint limits, self-collision, abnormal contact forces, or unstable hand motion are detected. |
| Articulation with floating gripper | Gripper jams into the object or cabinet frame | The handle trajectory is geometrically plausible but the gripper body collides with surrounding geometry. | Articulation progress stops, collision increases, or the gripper cannot maintain contact. |
| Insertion | Detected opening is too narrow or misaligned | The visual affordance is correct, but the insertion axis or clearance is insufficient. | Trajectory collides near the entrance or fails to reach insertion depth. |
| Hanging | Object falls after release | The support edge or hook does not create a gravity-stable contact. | After release, the object slides, rotates away, or drops under gravity. |
| Garment push / wash | Trajectory passes through IK singularity and moves erratically | The local end-effector path is valid, but the full arm has unstable or near-singular IK. | Full-manipulator validation shows large joint jumps, unstable end-effector motion, or IK failure. |
| Garment / deformable | Garment wraps around the arm or gripper | The trajectory sweeps through the cloth or pulls it across the robot body. | Arm-cloth collision, cloth entanglement, or task-outcome failure is detected. |
| Navigation | A* path is geometrically valid but dynamically infeasible | The grid path has sharp turns or narrow passages incompatible with the robot base. | DWB or base-aware rollout fails due to turning radius, velocity constraints, or clearance. |
| Navigation | Final base pose cannot support manipulation | The robot can navigate to the pose but cannot reach or see the manipulation target. | Manipulation reachability, visibility, or IK query at the final pose fails. |

Table 23: Common physics-aware augmentation strategies. Local perturbation explores nearby variations around validated anchors, while symmetry-aware augmentation expands candidates using symmetry priors from visual-language annotations.

| Skill / object type | Augmented quantity | Augmentation rule | Required recheck |
|--------------------------|---------------------------------------|---|--|
| Parallel-jaw grasp | Contact center | Jitter the grasp center within the graspable part mask or affordance region. | Target remains in region; gripper width valid; no collision; grasp remains stable. |
| Parallel-jaw grasp | Grasp rotation | Perturb in-plane rotation, approach yaw, or pre-grasp offset within a bounded range. | Approach path remains collision-free; contact is preserved; gravity robustness remains valid. |
| Dexterous grasp | Fingertip contacts | Perturb contact seeds on the same functional region and re-adjust hand configuration. | Contacts remain on region; hand joint limits, self-collision, and object stability are valid. |
| Dexterous grasp | Palm pose / joint offsets | Apply small palm-pose or finger-joint perturbations around a validated grasp. | No hand-object penetration; no self-collision; object remains constrained. |
| Garment fling | Grasp keypoints | Jitter sleeve endpoints, hem points, or bottom corners within keypoint confidence regions. | Bimanual reachability, visibility, cloth clearance, and fling outcome remain valid. |
| Garment fold | Fold/place target | Randomize sleeve place point, bottom-up fold point, or pants centerline target within a bounded region. | Fold target remains consistent with garment keypoints; no IK singularity or cloth-arm entanglement. |
| Articulation | Handle contact point | Jitter the contact point on the handle or movable-part mask. | Contact remains on handle; grasp or push remains feasible; no jamming. |
| Articulation | Pulling direction / way-point spacing | Perturb pull direction, way-point interval, and target joint range. | Joint progress remains valid; contact is maintained; gripper or hand does not collide with the frame. |
| Insertion | Alignment pose | Perturb entry pose, insertion depth, or axis orientation within tolerance. | Axis alignment, clearance, collision-free approach, and insertion depth remain valid. |
| Hanging | Support contact | Jitter the hook or support-edge contact point and release pose. | Gravity-stable support, clearance, and post-release hanging stability remain valid. |
| Navigation | Base pose | Perturb (x, y, ψ) around an interaction-ready base pose. | Traversability, clearance, path connectivity, DWB feasibility, and manipulation reachability remain valid. |
| Bottle / can | Rotationally symmetric grasp | Rotate validated grasp poses around the annotated symmetry axis. | Symmetry transform preserves object geometry; transformed grasp remains collision-free and stable. |
| Round knob | Rotationally symmetric contact | Rotate grasp or turning candidates around the knob axis. | Contact remains on knob; rotation direction and wrist feasibility remain valid. |
| Repeated cabinet handles | Translated or mirrored handle action | Copy grasp or articulation candidate to another functionally equivalent handle. | Target handle is equivalent; opening direction and collision context remain valid. |
| Garment structure | bilateral Mirrored sleeve fold | Mirror left-sleeve fold to right-sleeve fold across the garment centerline. | Mirrored keypoints are detected; fold direction and bimanual reachability remain valid. |