

Guiding Unified Dexterous Grasp Synthesis Across Modes and Scales via Learned Human Priors

Anonymous Author(s)
Affiliation
Address
email



Figure 1: **Human priors guide scalable dexterous grasp synthesis across modes and scales.** HUGS learns an object-conditioned human prior from a compact self-collected dataset to predict preferred contact modes and wrist initializations. Guided by this prior, force-closure-aware optimization synthesizes diverse and stable grasps ranging from two-finger pinches to bimanual grasps.

1 **Abstract:** Dexterous grasping across diverse object scales requires contact modes
 2 ranging from two-finger pinches to bimanual grasps. Existing dexterous grasp
 3 synthesis methods reduce the high-dimensional optimization space with manually
 4 designed expected contacts and initialization heuristics, which struggle to balance
 5 synthesis success rate and diversity. We present **HUGS**, a **H**uman-prior-guided
 6 framework for **U**nified dexterous **G**rasp **S**ynthesis across modes and scales. Instead
 7 of directly retargeting human demonstrations, HUGS learns an object-conditioned
 8 human prior that captures human grasp preferences and guides downstream force-
 9 closure-aware optimization. The prior is trained on a compact self-collected human
 10 grasp dataset with 1.8K grasps over 304 objects, providing broad coverage of
 11 object scales and contact modes. During synthesis, HUGS adaptively proposes
 12 contact modes and wrist initializations, substantially improving the balance between
 13 contact-mode coverage and synthesis success rate over heuristic-based methods.
 14 With HUGS, we synthesize 3.2M robotic grasps over 157K scenes, spanning object
 15 half-diagonal lengths from 2 cm to 30 cm and modes from two-finger to bimanual

16 grasps. Models trained on the synthesized dataset autonomously select appropriate
17 contact modes in the real world, enabling grasping from screws to large boxes.

18 **Keywords:** Dexterous Grasping, Human Priors, Grasp Synthesis

19 **1 Introduction**

20 Dexterous grasping is inherently multi-mode across object scales, from two-finger precision grasps
21 to coordinated bimanual grasps, and the same object may admit different modes under different
22 conditions. Training a generalizable grasp generation model for such diversity requires massive
23 datasets, yet large-scale real-world collection is impractical, making scalable grasp synthesis essential.

24 Dexterous grasp synthesis involves a vast search space due to the high degrees of freedom and
25 multi-contact nature. Existing methods often constrain optimization with manually predefined contact
26 modes, such as single-hand [1, 2] or dual-hand full-finger grasps [3, 4], and heuristic wrist pose
27 initializations. However, these coarse rules fail to exploit object-specific information and struggle to
28 balance synthesis efficiency and grasp diversity: loose heuristics lead to low-quality optimization and
29 inefficient synthesis, while restrictive heuristics severely limit grasp diversity.

30 Indeed, anthropomorphic robotic hands are designed for human-like manipulation [5], making human
31 grasps a natural source of knowledge for robots. Existing methods often retarget human grasps to
32 robotic hands [6, 7, 8], converting each demonstration into a robot grasp for the same or highly
33 similar object and thereby limiting scalable synthesis over large-scale, diverse object sets.

34 Our key insight is that we can use human demonstrations to learn generalized and reusable **object-**
35 **conditioned human priors**, instead of directly retargeting each human grasp to a robotic hand
36 in a one-to-one manner. By capturing human preferences for specific object geometry, the prior
37 guides robotic grasp optimization with better global initializations and optimization targets for
38 efficient convergence to higher-quality solutions. Such a **Human Prior + Robot Optimization**
39 paradigm offers following advantages: **1) Adaptivity:** compared with hand-crafted heuristics, the
40 learned human prior provides object-aware global guidance that constrains optimization to more
41 suitable regions of search space. **2) Scalability:** unlike direct retargeting methods, the learned prior
42 generalizes to unseen objects, enabling scalable synthesis of large and diverse robotic grasp datasets
43 from limited human demonstrations. and **3) Physical Plausibility:** force-aware optimization accounts
44 for robot-specific kinematics and task physics, while tolerating slight errors in coarse human prior.

45 In this work, we present **HUGS**, a Human-prior-guided framework for Unified dexterous Grasp
46 Synthesis across modes and scales. To learn a generalizable human prior from limited data, we
47 abstract human preference into contact modes and wrist poses, the two components most critical for
48 robotic grasp synthesis. Using a self-collected dataset with 304 objects and 1.8K grasps, we train
49 an object-conditioned generative model that predicts human-preferred grasp configurations. During
50 synthesis, the learned prior proposes object-adaptive contact modes and wrist initializations, enabling
51 scalable and efficient synthesis of 3.2M grasps over 157K scenes, spanning objects with half-diagonal
52 lengths from 2 cm to 30 cm and contact modes from two-finger to bimanual grasps, forming what is,
53 to our knowledge, the first large-scale dataset with multiple modes per object (Fig. 1). Models trained
54 on the dataset adaptively select suitable contact modes, enabling grasping from screws to large boxes
55 in the real world. Our contributions are highlighted as follows:

- 56 **1. Human Grasp Dataset Across Modes and Scales.** We collect and publicly release a compact yet
57 diverse human grasp dataset spanning a broad range of object scales and contact modes, providing
58 a strong foundation for learning generalizable human grasp priors.
- 59 **2. Human-Prior-Guided Unified Grasp Synthesis.** We propose a human-prior-guided framework
60 that learns object-conditioned priors over contact modes and wrist poses for unified dexterous
61 grasp synthesis across scales and modes, eliminating manually designed coarse heuristics.
- 62 **3. Large-Scale Synthesis and Evaluation.** We demonstrate scalable synthesis of diverse dexterous
63 grasp datasets and systematically evaluate the framework in terms of cross-scale and cross-mode
64 prior prediction, synthesis quality, efficiency, diversity, and real-world grasping potential.

65 2 Related Work

66 **Dexterous Grasp Synthesis.** Analytical grasp synthesis methods optimize hand configurations with
 67 physics-based objectives such as force closure. To reduce the search space of high-DoF multi-contact
 68 systems, existing methods typically rely on manually predefined contact regions and heuristic wrist
 69 initialization strategies. Most prior work focuses on single-hand full-finger grasps [9, 1, 2, 10] or
 70 bimanual grasps [3, 4]. Different contact modes can be synthesized by changing predefined contact
 71 regions, e.g., two-finger [11, 12] or three-finger grasps [13, 12]. However, these modes are typically
 72 manually tied to object-scale ranges, causing each object scale to correspond to only one grasp
 73 mode and limiting contact-mode diversity. In addition, most methods randomly sample wrist poses
 74 around the object’s surface [1, 2, 3]. For bimanual grasping, random sampling ignores inter-hand
 75 coordination, while strict symmetry constraints restrict diversity [3, 12]. Hand-centric methods
 76 [14, 15] improve diversity and efficiency for single-hand floating grasps, but may be less suitable for
 77 bimanual grasping and environment-constrained settings such as tabletop grasping.

78 **Human Data for Grasp Synthesis.** Human demonstrations are widely used for functional grasp
 79 synthesis through retargeting-based approaches that transfer human grasps to robotic hands. Some
 80 methods directly retarget each human grasp to a robotic grasp [7, 16], while others synthesize grasps
 81 for similar objects from a small set of human grasp templates [6, 17, 18, 8]. However, limited by
 82 the scale of human functional grasp datasets, these approaches typically support only limited object
 83 categories and rely on similarity to existing templates, rather than generalizing to arbitrary objects.
 84 Some works learn object-conditioned contact maps from human data [19, 20], but remain limited
 85 to single-hand full-finger grasps and require highly accurate human contact predictions that may
 86 not align with robotic hand kinematics. In contrast, our method uses only high-level human priors,
 87 namely contact modes and wrist poses, while retaining sufficient flexibility for force-closure grasp
 88 optimization compatible with robotic hand kinematics.

89 **Online Grasp Generation.** Synthetic grasp datasets are widely used to train grasp generation
 90 networks from partial observations such as images or point clouds for real-time deployment [21, 22,
 91 23, 13], including grasp generation in cluttered environments [24, 25, 26]. A few works directly learn
 92 grasping policies with reinforcement learning, bypassing explicit grasp pose planning [27, 28, 29, 30,
 93 31]; however, they often suffer from limited grasp diversity and unnatural finger motions.

94 3 Method

95 3.1 Problem Formulation

96 **Grasp Optimization Problem.** Let o denote the target object with surface \mathcal{S}_o . For hand $h \in \{L, R\}$,
 97 let $\mathbf{q}^h \in \mathbb{R}^{d_q}$ be the joint configuration and $\mathbf{T}^h = (\mathbf{R}^h, \mathbf{t}^h) \in SE(3)$ be the wrist pose, with bimanual
 98 states $\mathbf{Q} = (\mathbf{q}^L, \mathbf{q}^R)$ and $\mathbf{T} = (\mathbf{T}^L, \mathbf{T}^R)$. A grasp $g = (\mathbf{Q}, \mathbf{T}, c)$ uses contact configuration c to
 99 specify which hand regions contact the object, and grasp synthesis models $\pi(g | o) = \pi(\mathbf{Q}, \mathbf{T}, c | o)$.
 100 Active contact regions c induce contact points $\{\mathbf{p}_i\}_{i=1}^{N_c}$ with normals $\{\mathbf{n}_i\}_{i=1}^{N_c}$. Under point contact
 101 with Coulomb friction, each contact force satisfies $\mathbf{f}_i \in \mathcal{F}_i = \{\mathbf{f}_i | \|\mathbf{f}_i^t\| \leq \mu f_i^n, f_i^n \geq 0\}$ and
 102 induces wrench $\mathbf{w}_i = [\mathbf{f}_i; (\mathbf{p}_i - \mathbf{x}_{\text{com}}) \times \mathbf{f}_i]$, where \mathbf{x}_{com} is the object center of mass. The grasp
 103 wrench space is $\mathcal{G}(g) = \{\sum_{i=1}^{N_c} \mathbf{w}_i | \mathbf{f}_i \in \mathcal{F}_i\}$, and force closure holds when $0 \in \text{int}(\mathcal{G}(g))$. We
 104 formulate grasp optimization as minimizing force-closure energy $\Phi(g)$ under feasibility constraints:

$$\min_{g=(\mathbf{Q}, \mathbf{T}, c)} \Phi(g) \quad \text{s.t.} \quad \mathbf{Q} \in [\mathbf{Q}_{\min}, \mathbf{Q}_{\max}], \mathbf{p}_i \in \mathcal{S}_o \forall i, \text{CollisionFree}(g, o). \quad (1)$$

105 Here \mathbf{p}_i are induced by active regions in c ; the constraints enforce joint limits, object-surface contacts,
 106 and collision-free hand-object and hand-hand configurations.

107 **Overview of HUGS.** In grasp pose optimization, the contact configuration c determines the number
 108 and locations of contact points, making the problem hybrid discrete-continuous. Existing methods
 109 typically reduce the combinatorial search space with predefined contact regions. Meanwhile, the wrist
 110 pose \mathbf{T} is highly global, and poor initializations often trap local optimization in suboptimal minima.

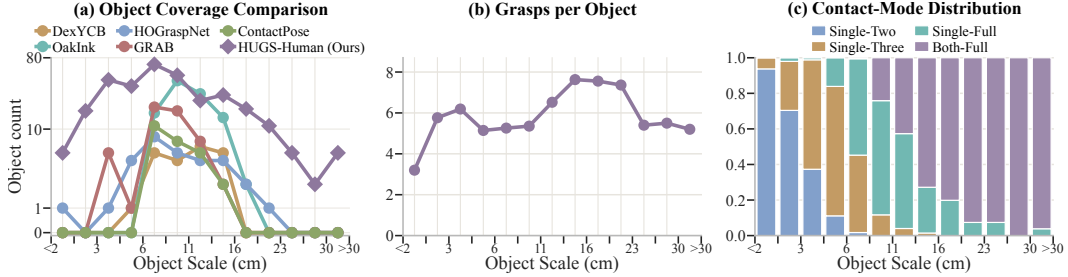


Figure 2: **Statistics of HUGS-Human Dataset.** (a) Object-count distribution, compared with existing human grasp datasets (OakInk: only real-world captured objects; log-scaled y-axis and power-scaled x-axis, with absolute-value labels). (b) Average number of HUGS-Human grasps per object. (c) Contact-mode distribution in HUGS-Human.

111 Consequently, optimization quality largely depends on initializing c and T , while optimizing the hand
 112 joint configuration Q is relatively straightforward given suitable c and T . Existing methods mainly
 113 use fixed coarse heuristics to initialize contact configurations and wrist poses, wasting optimization
 114 budget on implausible grasp modes while missing object-aware strategies. HUGS instead learns an
 115 object-conditioned human grasp prior to organize the search space before robot-specific optimization.

116 We adopt a compact discrete contact-mode representation with four coarse modes, selected as
 117 prevalent human grasping strategies that current dexterous robot hands can execute: 1) **Single-
 118 Two**, a single-hand two-finger grasp for small objects with limited contact area; 2) **Single-Three**, a
 119 single-hand three-finger grasp, the minimum for force closure under frictional point contacts [32]; 3)
 120 **Single-Full**, a denser single-hand full-finger grasp requiring larger accessible object surface area;
 121 and 4) **Both-Full**, a full-finger bimanual grasp for large objects requiring broader surface coverage.
 122 Based on this contact-mode abstraction, we formulate the grasp synthesis as

$$\pi(Q, T, c | o) = \pi(Q, T | T_0, c, o) \pi(T_0 | c, o) \pi(c | o), \quad (2)$$

123 where the learned prior predicts $\pi(c | o)$ and $\pi(T_0 | c, o)$ to propose plausible contact modes c and
 124 wrist initializations T_0 . Conditioned on these high-level proposals, the robot optimizer solves for Q
 125 and refines T under the force-closure and feasibility constraints in Eq. 1.

126 3.2 Compact Human Grasp Dataset

127 **Dataset Motivation.** Existing human grasp datasets mainly target hand-object reconstruction rather
 128 than synthesis guidance. Most emphasize single-hand grasps with limited scale coverage, such as
 129 DexYCB [33], OakInk [34], and HOGraspNet [35]; those containing bimanual grasps, such as GRAB
 130 [36] and ContactPose [37], include them only sparsely. They also do not explicitly consider mode
 131 and wrist-pose diversity, making them less suitable for learning priors over c and T_0 across scales.

132 **Self-Collected Dataset.** We introduce **HUGS-Human**, a compact human grasp dataset designed to
 133 learn the high-level prior used by HUGS. Participants grasp and lift tabletop objects from multiple
 134 supporting poses. For each posed object, we encourage diverse human-plausible contact modes and
 135 wrist poses while avoiding duplicates. The dataset contains 304 canonical objects, 525 posed objects,
 136 and 1.8k distinct human grasps. The collected grasps reveal three consistent patterns: 1) humans
 137 adapt contact modes and wrist poses according to object geometry and scale; 2) a single object
 138 can support multiple plausible contact modes, especially when its geometry is asymmetric; and 3)
 139 contact modes and feasible wrist poses are strongly coupled, as certain wrist poses only permit sparse
 140 fingertip contacts while others enable stable full-hand grasps on the same object. These patterns
 141 motivate our object-conditioned prior over c and T_0 . Each collected data sample includes the object
 142 mesh, object pose, hand wrist pose, MANO pose [38], and contact mode. The setup and annotation
 143 details are provided in the appendix. As shown in Fig. 2, HUGS-Human provides broad object-scale
 144 coverage compared with existing datasets, especially for very small and large objects, maintains
 145 multiple grasps per object across object scale bins, and captures the scale-dependent transition from
 146 sparse single-hand modes to full-hand and bimanual modes.

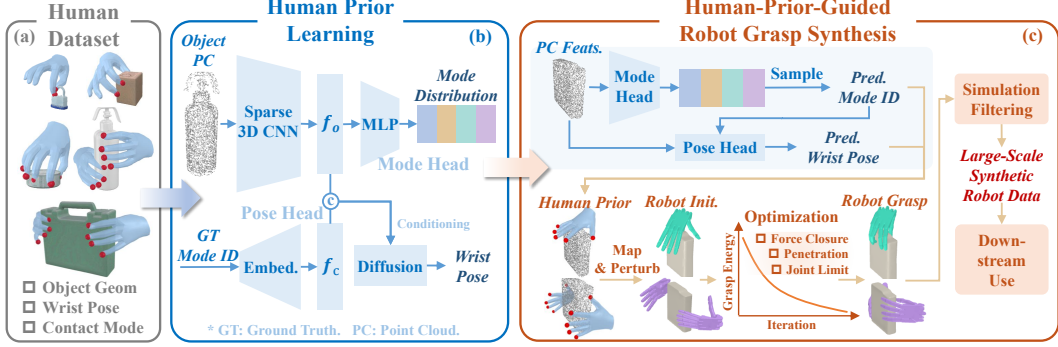


Figure 3: **Overview of HUGS.** HUGS learns an object-conditioned human prior from demonstrations with geometry, wrist poses, and contact modes. During synthesis, it predicts contact-mode distributions and wrist initializations for unseen objects, transfers them to robot hands, and optimizes force closure under feasibility constraints to synthesize large-scale robot grasps across modes and scales.

147 3.3 Object-Conditioned Human Grasp Prior

148 **Prior Formulation.** We learn object-conditioned distributions over contact modes and wrist poses,
 149 denoted by $\pi(c | o)$ and $\pi(\mathbf{T}_0 | c, o)$. Because human and robot wrist frames may differ, \mathbf{T}_0 uses the
 150 index metacarpophalangeal (MCP) position for translation and the dorsal-hand wrist orientation for
 151 rotation. We interpret these distributions as a **synthesis-budget prior**: for each object, they specify
 152 how much optimization effort to allocate to each contact mode and wrist-pose region. Implausible
 153 modes receive near-zero budget, while plausible modes share the budget according to human grasp
 154 diversity. This guides synthesis toward diverse plausible grasps with fewer optimization attempts.

155 **Architecture and Training.** As shown in Fig. 3(b), the prior takes a centered tabletop-frame object
 156 point cloud and encodes it into f_o using a MinkowskiEngine Sparse3DConv backbone [39]. A
 157 mode head predicts $\pi(c | o)$ with an MLP and softmax, supervised by the empirical contact-mode
 158 distribution of demonstrations for the same posed object using soft-label cross entropy. A mode-
 159 conditioned diffusion head models $\pi(\mathbf{T}_0 | c, o)$ by conditioning on $f_v = [f_o, f_c]$, where f_c is a
 160 learnable mode embedding, and is trained on individual data samples with contact-mode labels and
 161 human wrist poses. We use random vertical rotations and point-cloud noise for augmentation.

162 3.4 Human-Prior-Guided Grasp Synthesis

163 **Synthesis Overview.** During grasp synthesis, the learned human prior predicts $\pi(c | o)$ and
 164 $\pi(\mathbf{T}_0 | c, o)$. To account for the size differences between human and robot hands, we query the prior
 165 with an object point cloud rescaled by the human-to-robot hand scale ratio. Given sampled modes
 166 and wrist initializations, optimization produces robot grasps following $\pi(\mathbf{Q}, \mathbf{T} | \mathbf{T}_0, c, o)$. Together,
 167 they induce the final grasp distribution $\pi(\mathbf{Q}, \mathbf{T}, c | o)$ for cross-mode synthesis. Finally, we filter
 168 grasps in MuJoCo [40] and retain all validated grasps across contact modes for each object, rather
 169 than only the best one [12], to preserve dataset diversity for downstream applications.

170 **Prior-Guided Optimization Budget Allocation.** Before optimization, the learned prior predicts
 171 $\pi(c | o)$ and samples wrist initializations for each contact mode. For each mode, we draw $4K$
 172 candidates, rank them by approximate diffusion likelihood, keep the top $2K$, and apply farthest-pose
 173 selection to retain K diverse human-like initializations. The retained poses are mapped to the robot
 174 hand and slightly perturbed before optimization. Given a total budget of B optimization attempts, we
 175 allocate the mode budget as $B_c = \pi(c | o)B$ with a maximum per-mode cap. Each attempt samples
 176 one retained wrist initialization from its mode and open-hand joint angles.

177 **Grasp Optimization Formulation.** Given a contact mode c and wrist initialization \mathbf{T}_0 , we follow
 178 [2] and formulate grasp synthesis as a bi-level optimization, where the outer problem updates (\mathbf{Q}, \mathbf{T})
 179 and the inner problem evaluates force closure. We instantiate the force-closure energy in Eq. 1 as
 180 $\Phi(g) = \sum_{j=1}^s \Phi_j(g)$ over disturbance directions $\{t_j\}_{j=1}^s$, where each directional residual is the

181 optimal value of the following quadratic program (QP):

$$\Phi_j(g) \triangleq \min_{\mathbf{F}_j} \left\| \beta \mathbf{t}_j - \sum_{i=1}^{N_c} \mathbf{G}_i(g) \mathbf{f}_{j,i} \right\|^2 \quad \text{s.t. } \mathbf{f}_{j,i} \in \mathcal{F}_i, i = 1, \dots, N_c; \sum_{i=1}^{N_c} f_{j,i}^n \geq \gamma. \quad (3)$$

182 Here, $\mathbf{F}_j = \{\mathbf{f}_{j,i}\}_{i=1}^{N_c}$ are contact forces and $\mathbf{G}_i(g)$ is grasp matrix of contact point \mathbf{p}_i ; Φ_j measures
 183 resistance to disturbance \mathbf{t}_j . Unlike [12, 4], we use a weighted sum of global and per-hand force-
 184 closure terms for bimanual synthesis to promote both bimanual stability and individual hand quality.

185 4 Results

186 **Evaluation Goals.** We evaluate whether our object-geometry-aware human prior improves large-scale
 187 multi-mode synthesis through better contact mode allocation and wrist initialization, and supports
 188 downstream model distillation for real-world deployment across object scales.

189 **Implementation and Baselines.** We construct object scenes from DGN2k [2], which includes
 190 2.4k meshes. For each object, we sample multiple tabletop-stable poses and rescale it to 12 target
 191 sizes with AABB half-diagonal lengths from 2 cm to 30 cm, yielding 157k scenes spanning diverse
 192 geometries, poses, and scales. Each scene uses $B = 40$ optimization attempts in total, with at most
 193 $B_c = 20$ attempts per mode. Main results are reported on the Shadow Hand, with LEAP Hand [41]
 194 results in the appendix. Averaged over contact modes and object scales, HUGS optimizes 28.6 grasps
 195 per second on an RTX 4090. We compare HUGS against heuristic and ablated variants that differ in
 196 how contact modes and wrist-pose initializations are allocated. For the heuristic variants, scalar-scale
 197 contact-mode rules are derived from our human grasp dataset, and wrist poses are randomly initialized
 198 by convex-hull sampling with palm-facing poses. Further details are provided in the appendix.

- 199 • **Heur-Fix.** Uses the fixed *Single-Full* contact mode for every target object and scale, similar to [2].
- 200 • **Heur-Single.** Assigns one hand-designed contact mode to each object scale bin, similar to [12].
- 201 • **Heur-Multi.** Assigns multiple hand-designed contact modes to each object scale bin.
- 202 • **HUGS-Single.** Uses the HUGS wrist-pose prior, but allows only one mode per scale bin.
- 203 • **HUGS.** Our full method uses the human prior to allocate modes and initialize wrists.

204 4.1 Does Human Prior Predict Contact Modes Better than Scalar-Scale Rules?

205 We evaluate whether the geometry-aware human prior predicts
 206 contact-mode preferences better than rules based only on scalar
 207 object scale. HUGS-Human objects are split 8:2 into train and
 208 test sets; scalar-scale baselines map the object diagonal length
 209 to mode distributions using training-set statistics. Table 1
 210 reports KL divergence, soft precision, and soft recall on test
 211 set, with metric and baseline details in the appendix. The human prior achieves lower KL and higher
 212 precision and recall, suggesting that contact modes depend on object geometry beyond scalar scale.

Table 1. Comparison of contact-mode prediction on held-out objects.

Method	KL ↓	Prec. ↑	Rec. ↑
Scale Rules	0.612	0.805	0.892
Human Prior	0.300	0.873	0.964

213 4.2 Does Human-Prior Guidance Improve Efficient Multi-Mode Grasp Synthesis?

214 We next evaluate how human-prior guidance improves synthesis. Fig. 4(a) shows per-mode budgets
 215 (light background bars) and success counts (dark foreground bars), and Fig. 4(b) reports overall
 216 success across object scales. The heuristic baselines reveal the limits of scale-only rules. **Heur-Fix**
 217 works only near the scale suited to *Single-Full*, showing the need for contact-mode adaptation. **Heur-**
 218 **Single** is more robust, but one scale-dependent mode cannot cover multiple valid strategies within the
 219 same scale. **Heur-Multi** adds mode diversity, yet wastes attempts on object-specific mismatches, such
 220 as single-hand attempts on 13–23 cm objects. **HUGS** instead allocates budget by object geometry,
 221 smoothly shifting from *Single-Two* to *Single-Three*, *Single-Full*, and finally *Both-Full* as scale
 222 grows. It also captures object-specific exceptions within the same scale: at 13–19 cm, HUGS makes
 223 fewer *Single-Full* attempts but targets the few objects that admit single-hand grasps more accurately,

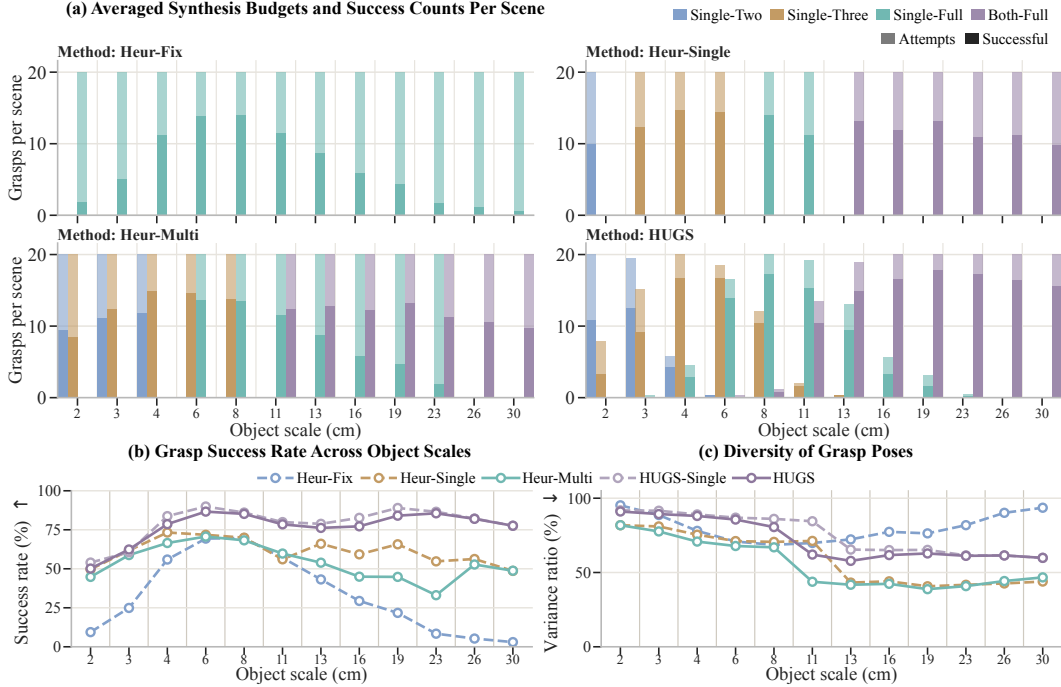


Figure 4: **Contact-mode allocation and synthesis success across object scales.** (a) Averaged synthesis budgets and success counts per scene. (b) Overall synthesis success rate across object scales. (c) Pose diversity measured by the explained-variance ratio of the first principal component.

224 yielding a much higher success rate than Heur-Multi. Thus, HUGS consistently outperforms all
 225 baselines. Fig. 1 visualizes diverse synthesized grasps, where many adjacent examples show the same
 226 object synthesized with different contact modes. Separately, HUGS-Single outperforms Heur-Single
 227 under the same single-mode setting, isolating the benefit of human-prior wrist initialization.

228 Fig. 4(c) compares grasp diversity using the explained-
 229 variance ratio of the first principal component [2, 14]. HUGS
 230 has lower pose dispersion because it concentrates on human-
 231 preferred wrist regions. This is expected: heuristic sampling
 232 increases diversity partly via unnatural reversed wrist poses
 233 (Fig. 5), which are difficult for humanoid robots to execute.

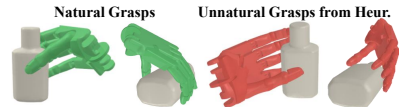


Figure 5: Heuristic wrist-pose sampling produces unnatural grasps.

234 4.3 Can Synthesized Grasps Supervise Online Grasp Generators?

235 We further test whether synthesized grasps can supervise learning-based grasp generation. We train
 236 identical lightweight object-conditioned generators on Heur-Multi and HUGS data, both of which
 237 contain multiple contact modes for the same object, using an 8:2 DGN2k object split and evaluating
 238 on held-out scenes across scales. The generator predicts binary contact-mode availability and
 239 corresponding grasp poses; architecture details are in the appendix. With HUGS data, contact-mode

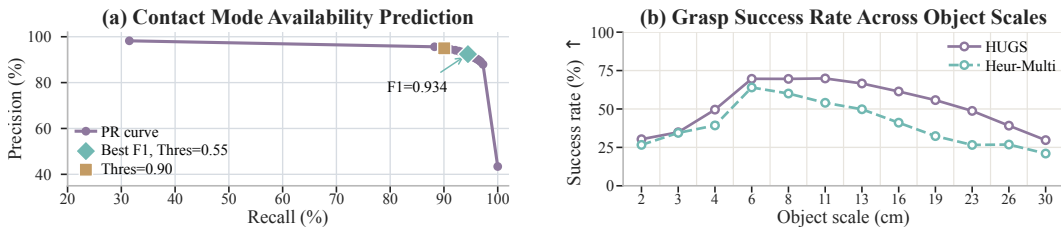


Figure 6: **Distilling from synthesized grasps.** (a) Precision-recall curve for contact-mode availability prediction. (b) Grasp success rate, comparing generators trained on HUGS and Heur-Multi data.

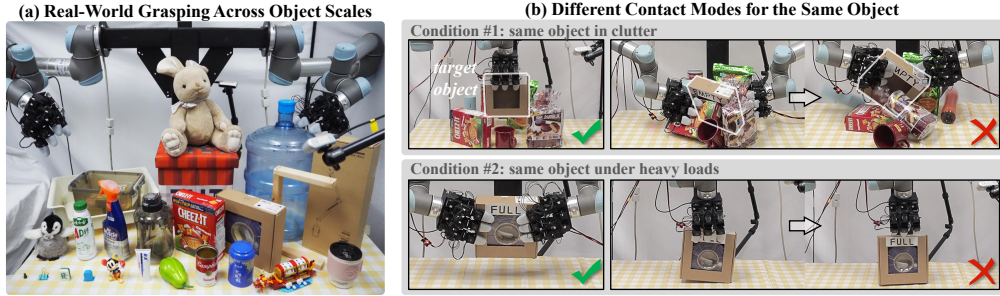


Figure 7: **Real-world grasping demonstrations on LEAP Hands.** (a) Cross-scale grasping over diverse objects. (b) Different mode selection under different deployment constraints.

240 availability is accurately predicted, reaching a best F1 score of 0.934 in Fig. 6(a). In simulation, the
 241 HUGS-trained generator achieves higher grasp success than the Heur-Multi-trained one (Fig. 6(b)),
 242 showing that HUGS produces easier-to-learn training data with better tolerance to generated-grasp
 243 errors. Success remains above 70% for medium objects (6–13 cm), but drops on very small and large
 244 objects, showing that cross-scale, cross-mode grasp generation remains challenging. Since our focus
 245 is grasp synthesis, we leave specialized generator design to future work.

246 4.4 Real-World Cross-Scale Grasping Demonstrations

247 We demonstrate real-world grasping with LEAP Hands using a generator trained from synthesized
 248 LEAP grasps. Across 24 objects, from a screw with a 0.5 cm half-diagonal to a bucket with a
 249 30 cm half-diagonal (Fig. 7(a)), the generator adaptively selects suitable contact modes and grasp
 250 poses. Fig. 7(b) shows why mode diversity matters even for the same object: in clutter, a top-down
 251 single-hand grasp remains feasible while side bimanual grasping is blocked; when the box is loaded,
 252 bimanual grasping becomes necessary to provide sufficient grasping force. Thus diverse generated
 253 grasps allow deployment-time constraints to select an appropriate grasp. More grasp cases are shown
 254 in Fig. 1 and the appendix; deployment details are also provided in the appendix. Note that the
 255 generator is still preliminary, and perception errors from calibration and depth reconstruction remain
 256 important bottlenecks for reliable grasping across object scales.

257 5 Limitations and Conclusion

258 **Limitations.** *Human data scale.* HUGS relies on a compact human dataset, so contact-mode and
 259 wrist-pose predictions may degrade on far out-of-distribution objects; scaling to large-scale egocentric
 260 interaction data is promising but requires further study of mode discovery and annotation. *Online*
 261 *grasp generation.* Our generator is preliminary and remains less reliable than the offline synthesis
 262 pipeline, especially for geometrically irregular objects, motivating research on stronger generative
 263 models and geometry representations. *Contact Diversity.* The four contact modes in HUGS capture
 264 dominant human grasp strategies but do not exhaustively cover the full spectrum of human grasping
 265 behaviors [42]. Extending to denser taxonomies or contact maps is needed for functional grasp
 266 synthesis beyond stable lifting. *Complete grasping systems.* Our real-world demonstrations remain
 267 sensitive to calibration, segmentation, and depth reconstruction errors, while robust deployment also
 268 requires end-to-end closed-loop execution, tactile adaptation, and collision-aware motion.

269 **Conclusion.** We presented HUGS, a human-prior-guided framework for unified dexterous grasp syn-
 270 thesis across contact modes and object scales. Instead of directly retargeting human demonstrations,
 271 HUGS learns an object-conditioned prior over contact modes and wrist poses to guide robot-specific
 272 optimization. Experiments show that this prior improves contact-mode prediction, synthesis success,
 273 and downstream grasp generation, while real-world demonstrations show adaptive mode selection
 274 and grasping from tiny screws to large containers. These results suggest that a human prior learned
 275 from compact human grasp data, when used as high-level guidance rather than direct imitation, can
 276 support scalable synthesis of diverse and executable dexterous robot grasps across modes and scales.

References

- [1] R. Wang, J. Zhang, J. Chen, Y. Xu, P. Li, T. Liu, and H. Wang. Dexgraspnet: A large-scale robotic dexterous grasp dataset for general objects based on simulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11359–11366. IEEE, 2023.
- [2] J. Chen, Y. Ke, and H. Wang. Bodex: Scalable and efficient robotic dexterous grasp synthesis using bilevel optimization. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 01–08. IEEE, 2025.
- [3] Y. Shao and C. Xiao. Bimanual grasp synthesis for dexterous robot hands. *IEEE Robotics and Automation Letters*, 9(12):11377–11384, 2024.
- [4] M. Lin, Y.-L. Wei, J. Chen, Y. Lin, S. Chen, J. Lyu, J. Chen, Y. Tang, H. Wang, and W.-S. Zheng. Bidexgrasp: Coordinated bimanual dexterous grasps across object geometries and sizes. *arXiv e-prints*, pages arXiv–2604, 2026.
- [5] A. Billard and D. Kragic. Trends and challenges in robot manipulation. *Science*, 364(6446): eaat8414, 2019.
- [6] W. Wei, P. Wang, S. Wang, Y. Luo, W. Li, D. Li, Y. Huang, and H. Duan. Learning human-like functional grasping for multifinger hands from few demonstrations. *IEEE Transactions on Robotics*, 40:3897–3916, 2024.
- [7] H. Chen, Y. Yao, Y. Ye, Z. Xu, H. Bharadhwaj, J. Wang, S. Tulsiani, Z. Erickson, and J. Ichnowski. Web2grasp: Learning functional grasps from web images of hand-object interactions. *arXiv preprint arXiv:2505.05517*, 2025.
- [8] S. Wang, Y. Yang, Y. Luo, D. Li, W. Wei, Y. Zhang, P. Hu, Y. Fu, H. Duan, J. Sun, et al. Scaleadfg: Affordance-based dexterous functional grasping via scalable dataset. *IEEE Robotics and Automation Letters*, 2025.
- [9] T. Liu, Z. Liu, Z. Jiao, Y. Zhu, and S.-C. Zhu. Synthesizing diverse and physically stable grasps with arbitrary hand structures using differentiable force closure estimator. *IEEE Robotics and Automation Letters*, 7(1):470–477, 2021.
- [10] R. Zurbrugg, A. Cramariuc, and M. Hutter. Graspqp: Differentiable optimization of force closure for diverse and robust dexterous grasping. In *Conference on Robot Learning*, 2025.
- [11] J. Ye, L. Wei, G. Jiang, C. Jing, X. Zou, and X. Wang. From power to precision: Learning fine-grained dexterity for multi-fingered robotic hands. *arXiv preprint arXiv:2511.13710*, 2025.
- [12] S. Yang, Y. Xie, Z. Liang, Y. Tian, J. Zeng, D. Lin, and J. Pang. Ultradexgrasp: Learning universal dexterous grasping for bimanual robots with synthetic data. *arXiv preprint arXiv:2603.05312*, 2026.
- [13] J. He, D. Li, X. Yu, Z. Qi, W. Zhang, J. Chen, Z. Zhang, Z. Zhang, L. Yi, and H. Wang. Dexvlg: Dexterous vision-language-grasp model at scale. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14248–14258, 2025.
- [14] J. Chen, Y. Ke, L. Peng, and H. Wang. Dexonomy: Synthesizing all dexterous grasp types in a grasp taxonomy. *arXiv preprint arXiv:2504.18829*, 2025.
- [15] Z.-H. Yin and P. Abbeel. Lightning grasp: High performance procedural grasp synthesis with contact fields. *arXiv preprint arXiv:2511.07418*, 2025.
- [16] Y.-L. Wei, J.-J. Jiang, C. Xing, X.-T. Tan, X.-M. Wu, H. Li, M. Cutkosky, and W.-S. Zheng. Grasp as you say: language-guided dexterous grasp generation. In *Proceedings of the 38th International Conference on Neural Information Processing Systems*, pages 46881–46907, 2024.

- 322 [17] C. Mao, H. Yuan, Z. Huang, C. Xu, K. Ma, and Z. Lu. Universal dexterous functional grasping
323 via demonstration-editing reinforcement learning. *arXiv preprint arXiv:2512.13380*, 2025.
- 324 [18] L. Huang, H. Zhang, Z. Wu, S. Christen, and J. Song. Fungrasp: Functional grasping for diverse
325 dexterous hands. *IEEE Robotics and Automation Letters*, 2025.
- 326 [19] Y. Ma, K. Chen, K. Zheng, and D. Qi. Contact map transfer with conditional diffusion model for
327 generalizable dexterous grasp generation. *Advances in Neural Information Processing Systems*,
328 38:79807–79836, 2026.
- 329 [20] Z. Wu, R. A. Potamias, X. Zhang, Z. Zhang, J. Deng, and S. Luo. Cedex: Cross-embodiment
330 dexterous grasp generation at scale from human-like contact representations. *arXiv preprint*
331 *arXiv:2509.24661*, 2025.
- 332 [21] Z. Wei, Z. Xu, J. Guo, Y. Hou, C. Gao, Z. Cai, J. Luo, and L. Shao. $\mathcal{D}(\mathcal{R}, \mathcal{O})$ grasp: A unified
333 representation of robot and object interaction for cross-embodiment dexterous grasping. In
334 *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4982–4988,
335 2025. doi:10.1109/ICRA55743.2025.11127754.
- 336 [22] Z. Weng, H. Lu, D. Kragic, and J. Lundell. Dexdiffuser: Generating dexterous grasps with
337 diffusion models. *IEEE Robotics and Automation Letters*, 9(12):11834–11840, 2024.
- 338 [23] Y. Xu, W. Wan, J. Zhang, H. Liu, Z. Shan, H. Shen, R. Wang, H. Geng, Y. Weng, J. Chen, et al.
339 Unidexgrasp: Universal robotic dexterous grasping via learning diverse proposal generation and
340 goal-conditioned policy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
341 *Pattern Recognition*, pages 4737–4746, 2023.
- 342 [24] J. Zhang, Z. Ma, T. Wu, Z. Chen, and H. Dong. Cadgrasp: Learning contact and collision aware
343 general dexterous grasping in cluttered scenes. *arXiv preprint arXiv:2601.15039*, 2026.
- 344 [25] Y. Zhong, X. Huang, R. Li, C. Zhang, Z. Chen, T. Guan, F. Zeng, K. N. Lui, Y. Ye, Y. Liang,
345 et al. Dexgraspvla: A vision-language-action framework towards general dexterous grasping. In
346 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pages 18836–18844,
347 2026.
- 348 [26] J. Zhang, H. Liu, D. Li, X. Yu, H. Geng, Y. Ding, J. Chen, and H. Wang. Dexgraspnet 2.0:
349 Learning generative dexterous grasping in large-scale synthetic cluttered scenes. In *8th Annual*
350 *Conference on Robot Learning*, 2024.
- 351 [27] H. Yuan, Z. Huang, Y. Wang, C. Mao, C. Xu, and Z. Lu. Demograsp: Universal dexterous
352 grasping from a single demonstration. *arXiv preprint arXiv:2509.22149*, 2025.
- 353 [28] H. Zhang, Z. Wu, L. Huang, S. Christen, and J. Song. Robustdexgrasp: Robust dexterous
354 grasping of general objects. *arXiv preprint arXiv:2504.05287*, 2025.
- 355 [29] Z. Chen, Q. Yan, Y. Chen, T. Wu, J. Zhang, Z. Ding, J. Li, Y. Yang, and H. Dong. Clutterdex-
356 grasp: A sim-to-real system for general dexterous grasping in cluttered scenes. In *Conference*
357 *on Robot Learning*, pages 885–905. PMLR, 2025.
- 358 [30] T. G. W. Lum, M. Matak, V. Makoviychuk, A. Handa, A. Allshire, T. Hermans, N. D. Ratliff,
359 and K. Van Wyk. Dextrah-g: Pixels-to-action dexterous arm-hand grasping with geometric
360 fabrics. In *Conference on Robot Learning*, pages 3182–3211. PMLR, 2025.
- 361 [31] R. Singh, A. Allshire, A. Handa, N. Ratliff, and K. Van Wyk. Dextrah-rgb: Visuomotor policies
362 to grasp anything with dexterous hands. *arXiv preprint arXiv:2412.01791*, 2024.
- 363 [32] K. M. Lynch and F. C. Park. *Modern robotics*. Cambridge University Press, 2017.

- 364 [33] Y.-W. Chao, W. Yang, Y. Xiang, P. Molchanov, A. Handa, J. Tremblay, Y. S. Narang, K. Van Wyk,
365 U. Iqbal, S. Birchfield, et al. Dexycb: A benchmark for capturing hand grasping of objects. In
366 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages
367 9044–9053, 2021.
- 368 [34] L. Yang, K. Li, X. Zhan, F. Wu, A. Xu, L. Liu, and C. Lu. Oakink: A large-scale knowl-
369 edge repository for understanding hand-object interaction. In *Proceedings of the IEEE/CVF*
370 *conference on computer vision and pattern recognition*, pages 20953–20962, 2022.
- 371 [35] W. Cho, J. Lee, M. Yi, M. Kim, T. Woo, D. Kim, T. Ha, H. Lee, J.-H. Ryu, W. Woo, et al. Dense
372 hand-object (ho) graspnet with full grasping taxonomy and dynamics. In *European Conference*
373 *on Computer Vision*, pages 284–303. Springer, 2024.
- 374 [36] O. Taheri, N. Ghorbani, M. J. Black, and D. Tzionas. Grab: A dataset of whole-body human
375 grasping of objects. In *European conference on computer vision*, pages 581–600. Springer,
376 2020.
- 377 [37] S. Brahmbhatt, C. Tang, C. D. Twigg, C. C. Kemp, and J. Hays. Contactpose: A dataset of
378 grasps with object contact and hand pose. In *European Conference on Computer Vision*, pages
379 361–378. Springer, 2020.
- 380 [38] J. Romero, D. Tzionas, and M. J. Black. Embodied hands: modeling and capturing hands and
381 bodies together. *ACM Transactions on Graphics (TOG)*, 36(6):1–17, 2017.
- 382 [39] C. Choy, J. Gwak, and S. Savarese. 4d spatio-temporal convnets: Minkowski convolutional
383 neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
384 *Recognition*, pages 3075–3084, 2019.
- 385 [40] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012*
386 *IEEE/RSJ international conference on intelligent robots and systems*, pages 5026–5033. IEEE,
387 2012.
- 388 [41] K. Shaw, A. Agarwal, and D. Pathak. Leap hand: Low-cost, efficient, and anthropomorphic
389 hand for robot learning. *arXiv preprint arXiv:2309.06440*, 2023.
- 390 [42] T. Feix, J. Romero, H.-B. Schmiebmayer, A. M. Dollar, and D. Kragic. The grasp taxonomy of
391 human grasp types. *IEEE Transactions on human-machine systems*, 46(1):66–77, 2015.
- 392 [43] T. Ren, S. Liu, A. Zeng, J. Lin, K. Li, H. Cao, J. Chen, X. Huang, Y. Chen, F. Yan, Z. Zeng,
393 H. Zhang, F. Li, J. Yang, H. Li, Q. Jiang, and L. Zhang. Grounded sam: Assembling open-world
394 models for diverse visual tasks, 2024.
- 395 [44] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland,
396 L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár,
397 and C. Feichtenhofer. Sam 2: Segment anything in images and videos, 2024. URL <https://arxiv.org/abs/2408.00714>.
398
- 399 [45] X. Chen, F.-J. Chu, P. Gleize, K. J. Liang, A. Sax, H. Tang, W. Wang, M. Guo, T. Hardin, X. Li,
400 et al. Sam 3d: 3dfy anything in images. *arXiv preprint arXiv:2511.16624*, 2025.
- 401 [46] B. Wen, W. Yang, J. Kautz, and S. Birchfield. Foundationpose: Unified 6d pose estimation and
402 tracking of novel objects. In *Proceedings of the IEEE/CVF conference on computer vision and*
403 *pattern recognition*, pages 17868–17879, 2024.
- 404 [47] R. A. Potamias, J. Zhang, J. Deng, and S. Zafeiriou. Wilor: End-to-end 3d hand localization
405 and reconstruction in-the-wild. In *Proceedings of the Computer Vision and Pattern Recognition*
406 *Conference*, pages 12242–12254, 2025.

- 407 [48] B. Sundaralingam, S. K. S. Hari, A. Fishman, C. Garrett, K. Van Wyk, V. Blukis, A. Millane,
408 H. Oleynikova, A. Handa, F. Ramos, and N. Ratliff. Curobo: Parallelized collision-free robot
409 motion generation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*,
410 pages 8112–8119. IEEE, 2023.
- 411 [49] B. Tan, C. Sun, X. Qin, H. Adai, Z. Fu, T. Zhou, H. Zhang, et al. Masked depth modeling for
412 spatial perception. *arXiv preprint arXiv:2601.17895*, 2026.

413 A Implementation Details

414 A.1 HUGS-Human Dataset Details

415 **Hardware Setup.** Raw demonstrations are recorded with four RealSense D435
416 cameras around the tabletop workspace, as shown in Fig. 8. All cameras are calibrated for both intrinsic and extrinsic
417 parameters, and the tabletop coordinate frame is also calibrated. We record RGB-D streams from all
418 four cameras. The depth images are aligned with the RGB images, and all RGB-D streams from the
419 four camera views are approximately synchronized and captured at a resolution of 1280×720 at 15
420 Hz. The raw recordings cover the complete manipulation process, including approaching the object,
lifting it, and placing it back on the table.

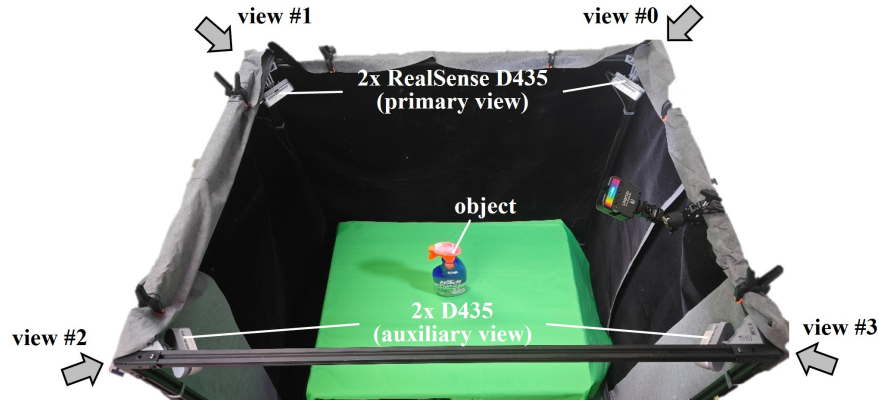


Figure 8: **Hardware setup for collecting the HUGS-Human dataset.** Four RealSense D435 cameras are used for hand-object mesh reconstruction. Camera views #0 and #1 serve as the primary views, while auxiliary views #2 and #3 are used to reduce single-view depth ambiguity for each hand.

421

422 **Annotation Details.** The dataset annotation includes hand-object mesh reconstruction and discrete
423 contact-mode annotation. We use camera views #0 and #1 as the primary views for reconstructing
424 the right and left hands, respectively, and auxiliary views #2 and #3 to reduce single-view depth
425 ambiguity for each hand. View #1 is used for object mesh reconstruction and pose tracking due to its
426 relatively low occlusion. The depth streams are used only to supervise object reconstruction. For
427 each raw recording, we use Grounded-SAM2 [43, 44] for object segmentation, SAM3D [45] for
428 mesh reconstruction, and FoundationPose [46] for pose estimation. The grasp frame is defined as
429 the frame immediately before clear object motion is detected. After detecting the grasp frame, we
430 use WiLoR [47] to estimate MANO hand poses [38] from the primary camera views #0 and #1. To
431 reduce the depth ambiguity inherent in single-view reconstruction, which can lead to hand-object
432 penetration or separation, we triangulate 3D hand keypoints from the primary and auxiliary views
433 using Direct Linear Transform (DLT). We then solve for the rigid transformation that best aligns
434 the WiLoR detections with the triangulated 3D keypoints. Finally, the discrete contact mode c
435 is manually annotated according to the number of hands involved in grasping and the number of
436 dominant contacting fingers, following the four modes defined in Sec. 3. For quality control, we
437 manually inspect all annotations and adjust them when necessary, such as correcting the selected
438 grasp frame. Each data sample includes the object mesh, object pose, hand wrist pose, MANO pose,
439 and contact mode. We will release both raw recordings and processed annotations.

440 **Dataset Visualization.** The 304 objects in HUGS-Human are collected from common office and
441 laboratory environments. We visualize the object-scale distribution, per-object grasp count, and
442 contact-mode distribution of HUGS-Human in Fig. 2. We further visualize representative samples
443 from the HUGS-Human dataset at different object scales in Fig. 9, showing both the extracted grasps
444 as reconstructed hand-object meshes and the raw RGB images from all four views. Finally, we

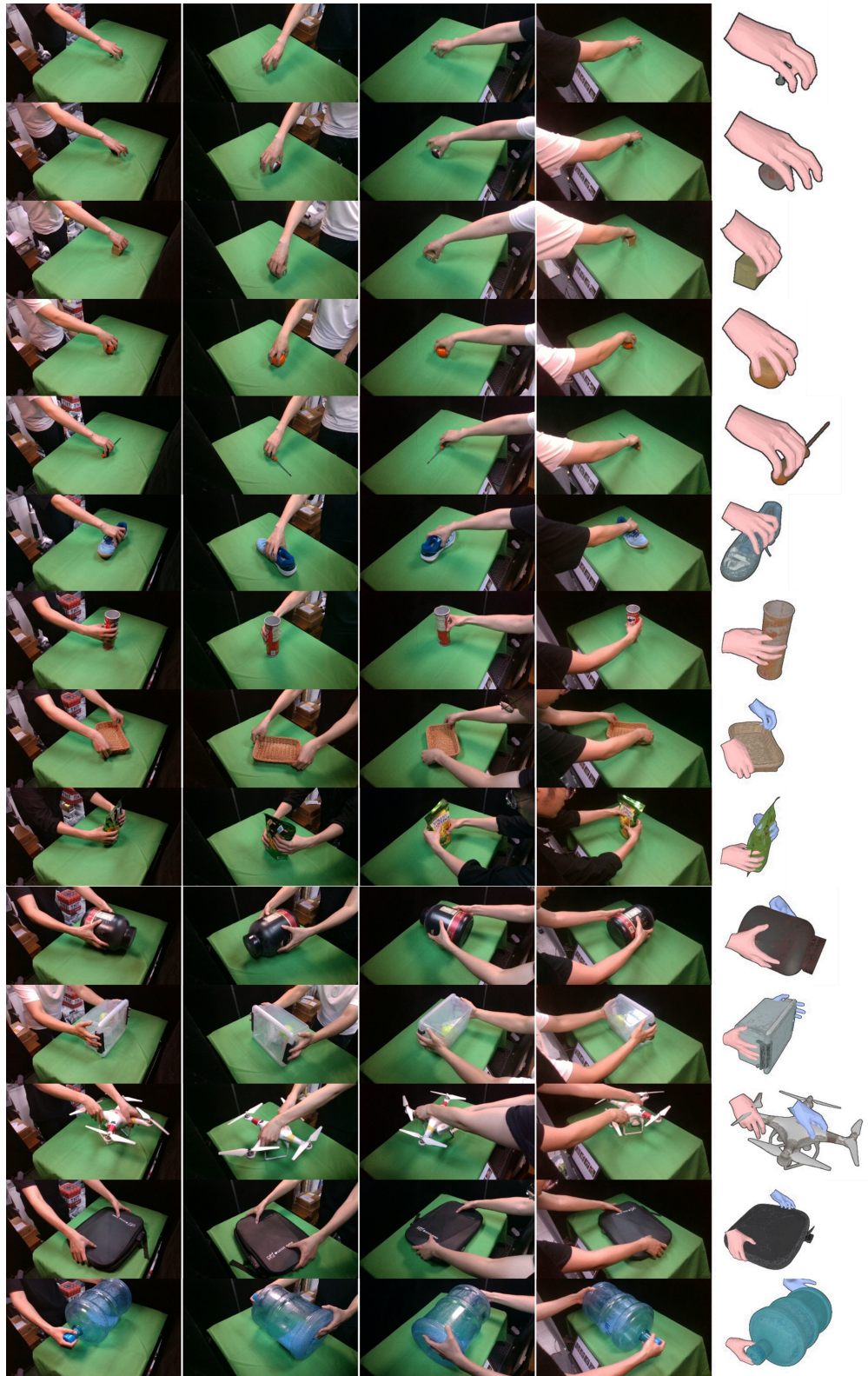
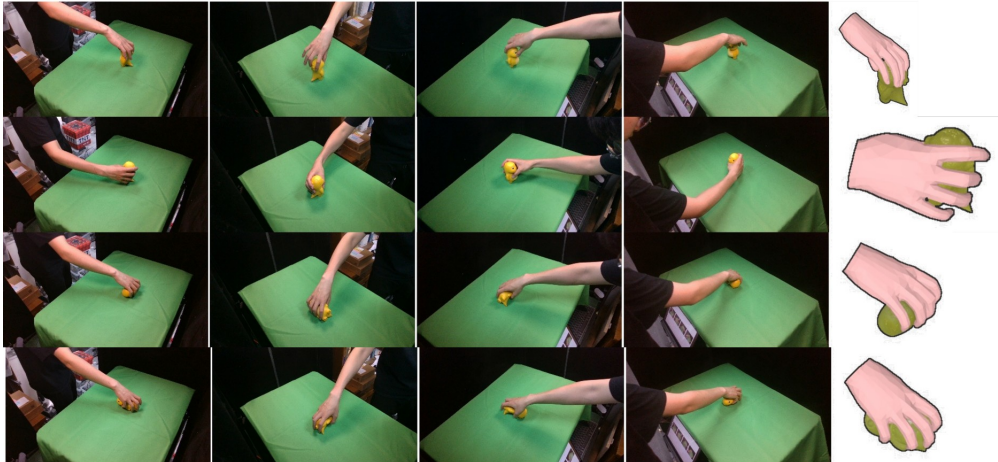


Figure 9: **Visualization of the HUGS-Human dataset.** Representative samples are shown with object scale increasing from top to bottom. Each sample includes RGB images from all four views, along with the rendered hand-object mesh reconstruction.

Object #1: the toy



Object #2: the watering can



Object #3: the vase

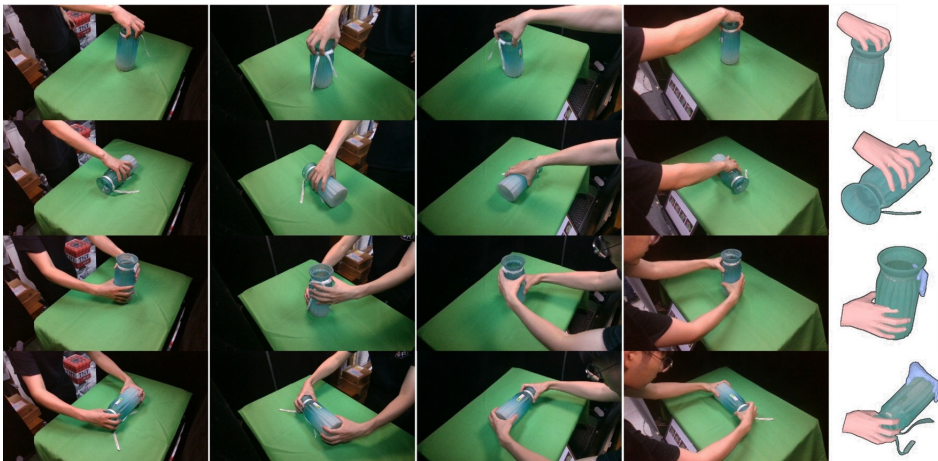


Figure 10: Visualization of representative cases of multiple grasps collected with a single object. For the smaller object #1, we collect single-handed grasps with varying numbers of fingers and object poses. For the larger objects #2 and #3, both single-handed and two-handed grasps are collected.

445 visualize representative cases of multiple grasps collected with a single object in Fig. 10, highlighting
 446 the object-level diversity of the HUGS-Human dataset.

447 **Rare Both-Three Contact Mode.** At the beginning of data collection, we also planned to include a
 448 *Both-Three* contact mode, where two hands are involved but each hand mainly uses sparse three-finger
 449 contacts. After completing the collection, however, we found that only a very small number of objects
 450 naturally support this mode, resulting in a negligible fraction of the dataset, as shown in Fig. 11.
 451 Incorporating *Both-Three* into the HUGS framework is technically feasible, since our contact-mode
 452 prior and optimization pipeline can be extended to additional discrete modes. Nevertheless, due to
 453 the limited number of valid demonstrations, evaluation results for this mode would be statistically
 454 unreliable. We therefore exclude *Both-Three* from the evaluation statistics and focus our analysis on
 455 the four dominant contact modes.

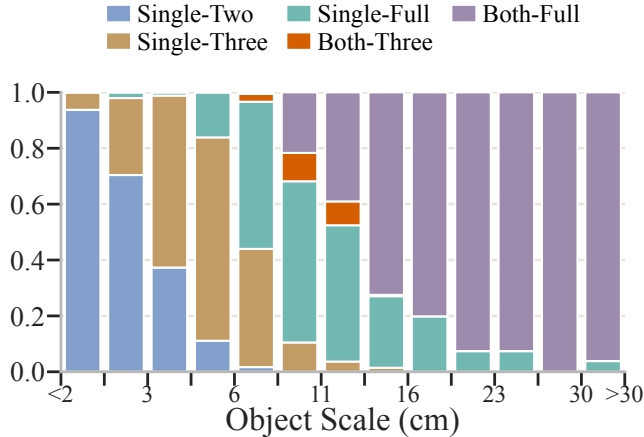


Figure 11: **Contact-mode distribution when including *Both-Three*.** The *Both-Three* mode accounts for only a negligible fraction of the collected demonstrations.

456 A.2 Human Prior Training Details

457 **Definition of Index MCP Frame.** To make the human wrist prior less sensitive to embodiment-
 458 specific wrist definitions, we define the predicted pose using an index MCP frame, as shown in
 459 Fig. 12. Its translation is the position of the index metacarpophalangeal (MCP) joint, and its rotation
 460 follows the dorsal-hand wrist orientation. This representation provides a stable hand-level pose that
 461 can be mapped to different robot hands during synthesis.

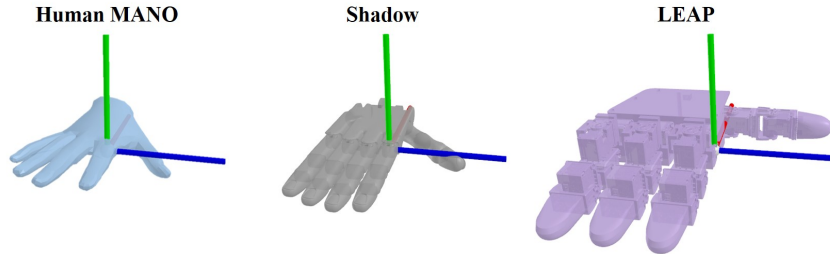


Figure 12: **Definition of the index MCP frame.** We use the index MCP position as the translation anchor and the dorsal-hand wrist orientation as the rotation reference for the human wrist prior.

462 **Architecture and Training.** In practice, we train the mode and wrist-pose heads separately because
 463 contact-mode prediction converges faster and overfits more easily than wrist-pose generation. Both
 464 models take the object point cloud as input. For each object sample, we uniformly sample 1024 points
 465 from the object point cloud and center the point cloud before feeding it to the network. The centered

466 point cloud is sparsely quantized with MinkowskiEngine using a voxel size of 0.005. The point-wise
 467 features are aggregated by mean pooling to obtain a 1024-dimensional global object feature. During
 468 training, we apply data augmentation including point-cloud centering, random rotations around the
 469 z axis, small rotations around the x and y axes with a maximum angle of 3° , and uniform scale
 470 augmentation in the range $[0.9, 1.1]$. For scale augmentation, valid hand translation targets are scaled
 471 accordingly. We also add Gaussian noise to point coordinates with standard deviation 0.001, clipped
 472 to three standard deviations, and apply point dropout with a ratio of 0.1. The contact-mode prior uses
 473 a two-layer MLP head. We supervise this branch with posed-object-level soft labels. For multiple
 474 human grasp records associated with the same posed object, we compute the occurrence frequency of
 475 each contact mode as the target distribution. It is trained with the soft-label cross-entropy loss

$$\mathcal{L}_{\text{mode}} = -\frac{1}{n_{\text{mode}}} \sum_i \sum_k q_{i,k} \log \text{softmax}(\mathbf{z}_i)_k, \quad (4)$$

476 where i indexes the samples included in the mode-supervision set, and k indexes the discrete mode
 477 classes. where q_i is the empirical contact-mode distribution computed from all human grasps of the
 478 same posed object in the dataset, and \mathbf{z}_i denotes the predicted logits. We train this branch from scratch
 479 for 100 iterations using AdamW with a batch size of 256. The learning rate is initialized to 10^{-3}
 480 and decayed to 10^{-4} with a cosine schedule. The pose prior is a contact-mode-conditioned diffusion
 481 model. We use a learnable 128-dimensional embedding for each contact mode and concatenate it with
 482 the 1024-dimensional global object feature to form a 1152-dimensional conditional feature. Training
 483 samples are uniformly drawn at random from all human grasps in HUGS-Human. The diffusion
 484 model generates a normalized 24-dimensional bimanual pose vector $[\mathbf{R}_r(9), \mathbf{p}_r(3), \mathbf{R}_l(9), \mathbf{p}_l(3)]$.
 485 The translation target \mathbf{p} represents the index-MCP position, and the rotation target \mathbf{R} uses the 9D
 486 representation of the wrist/palm rotation matrix. For right-only grasps, the left hand is filled with a
 487 fixed placeholder pose. During training, we apply RMS normalization to the pose vector. We use a
 488 cosine noise schedule with 128 training diffusion steps and optimize the denoiser with a velocity-
 489 prediction objective. The denoiser takes the noised 24D pose, the 1152-dimensional conditional
 490 feature, and the timestep embedding as inputs. It uses hidden layers with Mish activations and outputs
 491 a 24D velocity prediction. We train the denoiser with a SmoothL1 denoising objective between the
 492 predicted and target velocities. At inference time, we use a 10-step DDIM-style deterministic sampler
 493 conditioned on the object feature and contact-mode embedding. We train this branch from scratch for
 494 7500 iterations using AdamW with a batch size of 256. The learning rate is initialized to 10^{-3} and
 495 decayed to 10^{-4} with a cosine schedule.

496 A.3 Grasp Synthesis Details

497 **Contact Regions for Each Contact Mode.** Each contact mode specifies the active fingertip regions
 498 used by the robot hand during grasp optimization. For *Single-Two*, the active contacts are the thumb
 499 and index fingertips of one hand. For *Single-Three*, the active contacts are the thumb, index, and
 500 middle fingertips of one hand. For *Single-Full*, the active contacts are all fingertips of one hand. For
 501 *Both-Full*, the active contacts are all fingertips of both hands. During grasp optimization, different
 502 contact modes only change the expected fingertip contacts involved in the force-closure computation.

503 **Bimanual Optimization Implementation.** To adapt bimanual grasp optimization to the frameworks
 504 of [2, 48], we model the two hands as a single composite robot by adding dummy joints. The dummy
 505 joints include three translational joints and three rotational joints, which parameterize the floating
 506 wrist pose of each hand. In practice, very large objects may lie near or beyond the boundary of
 507 the human grasp dataset distribution, where training samples are relatively sparse. As a result, the
 508 predicted global hand pose can be slightly less accurate, leading to noticeable hand-object penetration
 509 in some cases. To make the initial states more suitable for subsequent optimization, we offset each
 510 hand by 10 cm along its palm-outward direction when constructing the initial robot wrist poses, as
 511 shown in Fig. 21. After sampling wrist poses from the human prior, we first solve batched inverse
 512 kinematics to obtain the corresponding dummy joint angles. These IK solutions are then used as the
 513 initialization for the subsequent grasp optimization. Unlike the default single-hand configuration
 514 in [2], bimanual optimization first optimizes the in-contact grasp pose, and then optimizes the

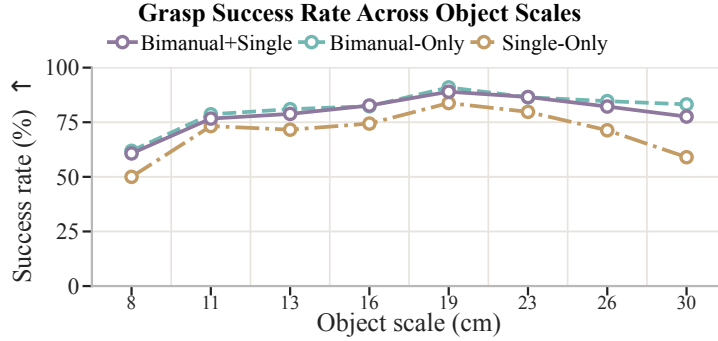


Figure 13: **Ablation of the bimanual force-closure objective.** We compare using only the per-hand force-closure terms (*Single-Only*), only the global bimanual force-closure term (*Bimanual-Only*), and their combination (*Bimanual+Single*). Without considering bimanual force closure, the grasp success rate drops significantly.

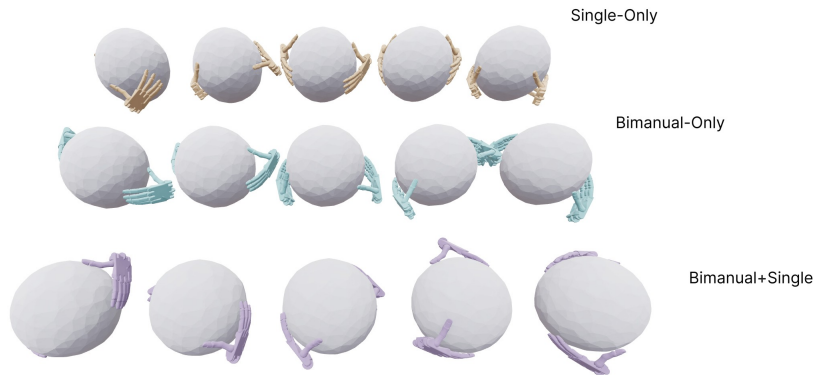


Figure 14: **Qualitative examples of bimanual force-closure ablation.** Using only the global bimanual term (*Bimanual-Only*) can produce degenerate contacts, where each hand touches the object only weakly or at fingertip extremities. Combining global and per-hand terms (*Bimanual+Single*) yields more coordinated bimanual grasps while maintaining better individual hand contact quality.

515 non-contact pregrasp pose in a final stage by reducing the hand-to-object point distance threshold
 516 by 1 cm. This ordering allows the grasp pose to be closer to the object center. In addition, we
 517 assign a larger step size to the dummy joint (base pose) variables, which decays over optimization
 518 iterations, allowing the wrist poses to adjust more aggressively in early iterations while converging
 519 stably. For the force-closure objective, we use a sum of a global bimanual term and per-hand terms:
 520 $\Phi_{bi}(g) = \Phi_{global}(g) + \Phi_{left}(g) + \Phi_{right}(g)$, where Φ_{global} evaluates force closure over contacts from
 521 both hands jointly, and Φ_{left} , Φ_{right} evaluate each hand independently. Fig. 13 shows that omitting
 522 the global bimanual force-closure term substantially reduces grasp success, indicating that per-hand
 523 force closure alone is insufficient for stable bimanual grasping. Although *Bimanual-Only* achieves a
 524 success rate similar to the combined objective, Fig. 14 shows that using only the global bimanual
 525 term can produce degenerate contacts, where each hand touches the object only weakly or at fingertip
 526 extremities. Therefore, the global term is needed for overall bimanual stability, while the per-hand
 527 terms improve the contact quality of each individual hand.

528 **Human-to-Robot Scale Ratio.** Before querying the human prior, we rescale the object point cloud
 529 according to the human-to-robot hand scale ratio. The motivation is to match the relative object size
 530 perceived by the robot hand to that in the human demonstrations, since different hand embodiments
 531 have different physical sizes as illustrated in Fig. 12. For example, the LEAP Hand is larger than a
 532 human hand, so an object grasped by a human appears relatively smaller to LEAP. Therefore, when

533 querying the human prior for LEAP synthesis, we downscale the object point cloud before prediction.
534 We use a scale ratio of 1.0 for the Shadow Hand and 1.4 for the LEAP Hand.

535 **MuJoCo Simulation.** We filter optimized grasps in MuJoCo simulation by lifting the object and
536 checking three criteria for success. First, the object must be lifted to a height of 0.1 m. Second, the
537 object pose must remain stable throughout the lift, with position deviation below 5 cm and rotation
538 deviation below 15°. Third, the grasp must be collision-free. A grasp is disqualified if the hand is
539 in self-collision or intersects the object at the start of the simulation, or if self-collision occurs at
540 any point during the grasping process. For physical simulation across object scales, we assign a
541 scale-dependent effective density to each object. To avoid cubic mass growth under geometric scaling,
542 we use $d(s) = d_0(s/s_0)^{-2}$, where $d_0 = 700 \text{ kg/m}^3$ and $s_0 = 0.06 \text{ m}$. This formulation makes the
543 resulting object mass scale approximately linearly with s . The resulting median object masses are
544 0.032 kg at scale 0.02, 0.096 kg at scale 0.06, 0.160 kg at scale 0.1, and 0.339 kg at scale 0.3.

545 A.4 Grasp Generator Details

546 **Architecture.** The robot grasp prior uses the same sparse point-cloud encoder as the human prior. The
547 input partial object point cloud is fused from three camera views. We sample 1024 points from the
548 partial object point cloud, apply MinkowskiEngine sparse quantization, and feed the quantized point
549 cloud into MinkUNet. Mean pooling over point-wise features produces a 1024-dimensional object
550 feature. The model contains a contact-mode availability head and a mode-conditioned robot pose
551 head. The availability head predicts independent binary availability scores for the four contact modes.
552 Its supervision is computed from the contact-mode records of the same scene in the synthetic robot
553 dataset, allowing multiple contact modes to be available simultaneously. The pose-head architecture
554 follows [14]. It conditions on the concatenation of the object feature and a learnable 128-dimensional
555 contact-mode embedding. It first uses a conditional diffusion model to generate the final grasp-frame
556 24-dimensional bimanual wrist pose $[\mathbf{R}_r, \mathbf{p}_r, \mathbf{R}_l, \mathbf{p}_l]$, where rotations use the 9D rotation-matrix
557 representation. The diffusion model uses RMS-normalized poses, 128 training diffusion steps, a
558 cosine noise schedule, and a velocity-prediction objective, and is trained with a SmoothL1 denoising
559 loss. To recover the full robot trajectory, we further use two MLPs for the left and right hands.
560 Conditioned on the final wrist pose, these MLPs regress wrist translations and rotations for the
561 pregrasp and grasp stages, as well as finger joint positions for the pregrasp, grasp, and squeeze stages.

562 **Training.** The overall training objective consists of five terms: contact-mode availability BCE
563 loss, final-wrist diffusion denoising loss, trajectory translation loss, trajectory rotation loss, and
564 joint-position loss. We use equal weights for all five terms in the current configuration and train the
565 entire network jointly in a single stage for 50000 iterations. We use AdamW with a batch size of 256.
566 The learning rate is initialized to 10^{-3} and decayed to 10^{-4} with a cosine schedule.

567 **Sampling.** At inference time, the availability head first predicts independent contact-mode scores
568 for the input partial point cloud. We select feasible contact modes using a score threshold of 0.9.
569 According to the precision-recall curve in Fig. 6(a), this threshold yields approximately 95% precision
570 and 90% recall for the Shadow Hand. For each selected contact mode, we sample the final wrist
571 pose with the DDIM deterministic sampler conditioned on the object feature and contact-mode
572 embedding, followed by RMS de-normalization. The trajectory MLPs then predict the pregrasp and
573 grasp wrist poses, as well as the pregrasp, grasp, and squeeze joint positions, conditioned on the
574 sampled final wrist pose. For each feasible contact mode, we sample 100 grasp poses and select
575 the top 10 candidates according to the approximate diffusion probability for simulation evaluation.
576 The resulting candidates may be further filtered by reachability, collision checking, and task-specific
577 deployment constraints for downstream applications.

578 A.5 Real-World Deployment Details

579 **Deployment Pipeline.** The real-world hardware setup is shown in Fig. 15. We deploy the system
580 on a dual-arm UR5 platform with two LEAP Hands equipped with customized fingertips. During
581 deployment, three RealSense D435 cameras capture RGB-D observations. All cameras are calibrated

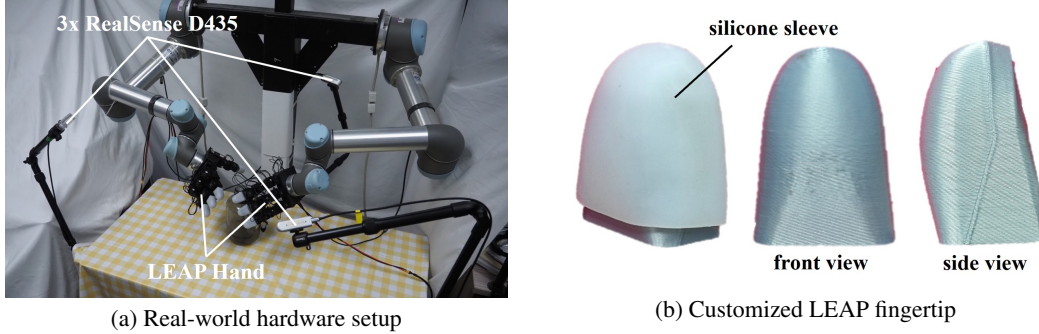


Figure 15: **Real-world hardware setup and customized LEAP fingertip.** Our hardware setup consists of three RealSense D435 cameras arranged approximately 120° apart for object point-cloud reconstruction. Compared with the original LEAP fingertip, our customized fingertip has a smoother geometry and is covered with a silicone sleeve to improve contact stability.

582 for inter-camera extrinsics, and we further perform hand-eye calibration between the robot and the
 583 primary camera. RGB images are segmented with Grounded-SAM2 [43, 44], depth maps are refined
 584 with Lingbot-depth [49], and the fused object point cloud is fed to the generator. Collision-free
 585 approach motions are planned using cuRobo [48]. For real-world execution, we sample 25 grasps for
 586 each feasible contact mode and manually select three candidates for robot execution. Snapshots of
 587 the executed real-world grasps on objects #1–#23 are shown in Fig. 22 and Fig. 23, while object #24
 588 is used to demonstrate the significance of contact-mode diversity (Fig. 7 (b)). As discussed in Sec. 5,
 589 real-world generation is affected by calibration errors, depth reconstruction noise, segmentation
 590 quality, and other deployment-specific factors. Moreover, our online generator is still preliminary
 591 and remains less reliable than the offline synthesis pipeline. Therefore, the real-world experiments
 592 are intended to demonstrate the deployment potential of HUGS in physical scenes, rather than to
 593 provide a quantitative evaluation of grasp success rate. The typical failure modes are analyzed in
 594 Appendix D.3. Building a fully autonomous and robust real-world grasping system remains an
 595 important direction for future work.

596 B Evaluation Details

597 B.1 Metrics Details

598 **Human Contact-Mode Distribution Prediction.** This section provides the metric definitions for the
 599 contact-mode distribution prediction experiment in Sec. 4. Given a posed object, the task is to predict
 600 the empirical distribution of human-preferred contact modes rather than a single contact-mode label.
 601 Let \mathbf{q} denote the empirical contact-mode distribution and \mathbf{p} denote the predicted distribution over
 602 contact modes. We measure the discrepancy between the two distributions using

$$D_{\text{KL}}(\mathbf{q}||\mathbf{p}) = \sum_t q_t \log \frac{q_t}{p_t}. \quad (5)$$

603 We also report soft precision,

$$\text{SoftPrec}(\mathbf{p}, \mathbf{q}) = \sum_{t \in \mathcal{P}(\mathbf{q})} p_t, \quad \mathcal{P}(\mathbf{q}) = \{t \mid q_t > 0\}, \quad (6)$$

604 which measures how much predicted probability mass falls on ground-truth positive contact modes.
 605 Finally, we report soft recall,

$$\text{SoftRec}(\mathbf{p}, \mathbf{q}) = \sum_{t \in \text{Top}_{|\mathcal{P}(\mathbf{q})|}(\mathbf{p})} q_t, \quad (7)$$

606 which measures how much ground-truth probability mass is covered by the top predicted modes, with
 607 the number of selected modes set to $|\mathcal{P}(\mathbf{q})|$, i.e., the number of ground-truth positive contact modes.
 608 Lower KL and higher soft precision/recall indicate better contact-mode prediction.

609 **Synthesis Success Rate.** For a given evaluation group, the synthesis success rate is defined as the
610 fraction of optimization attempts that pass physical validation in simulation, i.e., $\text{SuccessRate} =$
611 $N_{\text{success}}/N_{\text{attempt}}$, where N_{attempt} is the number of grasp optimization attempts and N_{success} is the
612 number of physically successful grasps.

613 **Success Counts per Scene.** For each object scene, the success count is the number of successful
614 grasps obtained from all optimization attempts allocated to that scene. When reporting success
615 counts per scene across an object-scale bin or a method, we average this count over all scenes in the
616 corresponding group.

617 **Pose Diversity.** We measure grasp-pose diversity using the explained-variance ratio of the first
618 principal component, following the PCA-based diversity metric in prior work [2, 14]. For each
619 synthesized grasp, we use the optimized wrist pose and hand joint angles as the pose feature,
620 including wrist translation, wrist rotation, and finger joint positions. For each object scene, we
621 perform PCA over all validated grasp poses of that scene and record the explained-variance ratio
622 of the first principal component. We then average this value over all scenes in the corresponding
623 evaluation group. This differs from [2, 14], which compute PCA after pooling grasps across scenes.
624 We adopt the scene-wise formulation because our focus is the diversity of valid grasp poses for the
625 same scene. A lower first-component explained-variance ratio indicates that successful grasps are
626 less dominated by a single pose direction and are therefore less concentrated along one major mode.
627 In our analysis, however, large dispersion can also arise from unnatural wrist poses produced by
628 heuristic sampling; thus the diversity metric should be interpreted together with grasp naturalness
629 and qualitative examples.

630 **Robot Contact-Mode Availability Prediction.** For the robot grasp generator, contact-mode
631 availability is evaluated as a multi-label binary prediction problem. For each scene and contact
632 mode, the ground-truth label is positive if the synthesized dataset contains at least one validated
633 grasp of that mode for the scene, and negative otherwise. Given predicted availability scores, we
634 sweep the decision threshold to compute precision and recall: $\text{Precision} = \text{TP}/(\text{TP} + \text{FP})$ and
635 $\text{Recall} = \text{TP}/(\text{TP} + \text{FN})$. The F1 score is defined as $\text{F1} = 2 \text{Precision Recall}/(\text{Precision} + \text{Recall})$.
636 We report the precision-recall curve and the best F1 score over thresholds.

637 B.2 Baseline Details

638 **Scalar-Scale Rules.** The scale-rule baseline is derived from statistics of HUGS-Human. Here, object
639 scale refers to the half-diagonal length of the object’s axis-aligned bounding box (AABB). We first
640 group posed objects in the training split into scale bins according to their object size. For each
641 scale bin, we aggregate all human grasps associated with posed objects in that bin and compute the
642 empirical contact-mode distribution. At test time, a posed object is assigned to the corresponding
643 scale bin based only on its scalar object size, and the precomputed bin-level contact-mode distribution
644 is used as the prediction. This baseline therefore captures scale-dependent mode preferences but
645 ignores detailed object geometry. In the contact-mode distribution prediction experiment in Sec. 4,
646 we evaluate the continuous contact-mode distributions directly. For grasp synthesis, however, the
647 scalar-scale rules are used only to decide whether each contact mode should be considered. All
648 contact modes selected by the rule are assigned the same optimization budget.

649 **Heuristic Wrist-Pose Sampling.** For heuristic baselines, we use convex-hull wrist sampling to
650 initialize optimization. For each object scene, we first construct an expanded convex hull of the
651 object. After applying the object scale in the scene, we expand the hull along the hull-vertex normals.
652 We then uniformly sample candidate points on the expanded hull surface and transform both the
653 sampled points and face normals into the world frame. Each candidate point defines a palm seed. The
654 translation is set to the sampled point, and the approach axis is set opposite to the outward normal so
655 that the palm faces the object. Candidate seeds are filtered by collision-free checks. For bimanual
656 grasping, we first randomly sample a right-hand surface candidate. We then compute the mean xy
657 position of all valid candidates in the current scene as the pairing center, mirror the right-hand xy
658 position about this center, and keep the right-hand z coordinate unchanged. In the left-hand candidate

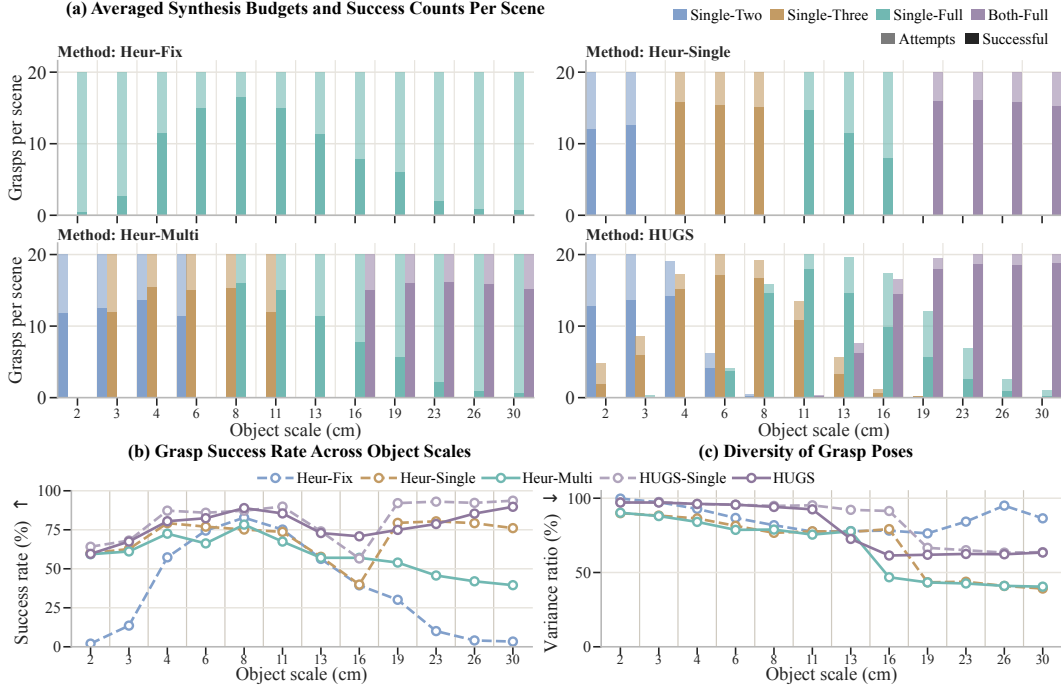


Figure 16: **Contact-mode allocation and synthesis success on the LEAP Hand.** (a) Averaged synthesis budgets and success counts per scene. (b) Overall synthesis success rate across object scales. (c) Pose diversity measured by the explained-variance ratio of the first principal component.

659 set, we select the real surface candidate whose translation is closest to this mirrored target. This
 660 procedure only pairs translations; the left-hand rotation is still constructed from the surface normal of
 661 the selected left-hand candidate, rather than by mirroring the right-hand orientation. Finally, we add
 662 translational and rotational perturbations to the sampled poses and convert them into initial poses.

663 C Additional Results

664 C.1 LEAP Hand Simulation Results

665 We provide additional LEAP Hand simulation results using the same synthesis protocol as the Shadow
 666 Hand experiments, except that we use a human-to-robot scale ratio of 1.4 when querying the human
 667 prior. Since the LEAP Hand is larger than a human hand, an object grasped by a human appears
 668 relatively smaller to LEAP; therefore, we downscale the object point cloud before prediction.

669 **LEAP Synthesis Results.** The LEAP Hand results in Fig. 16 show trends consistent with the Shadow
 670 Hand results in Sec. 4. Across object scales, HUGS adaptively reallocates synthesis budget across
 671 contact modes instead of relying on fixed scalar-scale rules. This object-conditioned allocation
 672 improves synthesis success over heuristic baselines while preserving multi-mode grasp coverage.
 673 The diversity trend is also consistent with the Shadow Hand experiments: HUGS concentrates
 674 optimization around human-preferred wrist regions, whereas heuristic sampling can introduce larger
 675 pose dispersion through less natural wrist initializations. These results suggest that the learned human
 676 prior is not tied to a single robot hand and can transfer to a different dexterous hand, even when its
 677 physical size differs from that of a human hand.

678 **Qualitative Synthetic LEAP Grasps.** Fig. 17 shows qualitative examples of HUGS-synthesized
 679 grasps on the LEAP Hand. The examples span different object scales and contact modes, from sparse
 680 single-hand grasps on small objects to full-hand and bimanual grasps on larger objects. They also
 681 include cases where different contact modes are synthesized for the same object.



Figure 17: Synthesized LEAP Hand grasps across object scales and contact modes.

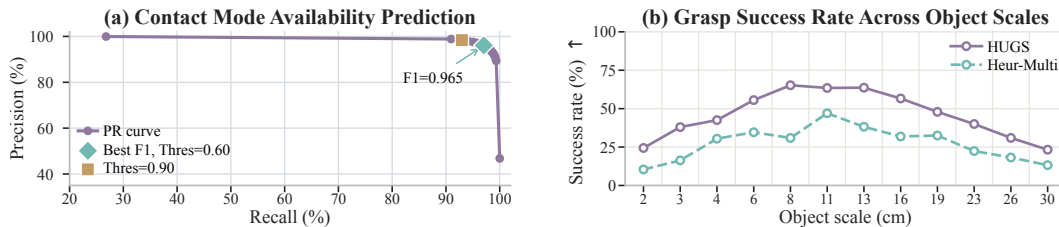


Figure 18: **Distilling from synthesized LEAP grasps.** (a) Precision-recall curve for contact-mode availability prediction. (b) Grasp success rate, comparing generators trained on HUGS and Heur-Multi LEAP grasp data.

682 **LEAP Distillation Results.** Fig. 18 further evaluates distillation from synthesized LEAP grasps,
 683 following the same protocol as the Shadow Hand generator experiment in Fig. 6. The results are
 684 consistent with the Shadow Hand setting: the generator trained on HUGS data predicts contact-mode
 685 availability and produces more successful grasps than the generator trained on heuristic data. This
 686 indicates that the synthesized HUGS grasps remain useful for supervising online grasp generators
 687 after transferring the synthesis pipeline to the LEAP Hand.

688 C.2 More Visualization

689 **Contact-Mode Prior Visualization.** Fig. 19 visualizes contact-mode distributions predicted by the
 690 learned human prior for diverse object point clouds. The examples show that the prior adapts its mode
 691 probabilities according to object geometry and scale, assigning high probability to sparse single-hand
 692 modes for small objects and to full-hand or bimanual modes for larger objects.

693 **Wrist-Pose Prior Visualization.** Fig. 20 shows mode-conditioned wrist-pose samples from the
 694 human prior. For the same object, different contact modes can induce different preferred approach
 695 regions. In Case 1, the object is large and wide. Bimanual grasps can approach either from the side
 696 or from above because the two hands can jointly cover the object, whereas single-hand side grasps
 697 would place the contact region far from the object’s center of mass. Therefore, the single-hand prior
 698 favors top-down grasps whose contact region is closer to the gravity line through the center of mass.
 699 In Case 2, the object is approximately cylindrical with height larger than its radius. From a top-down
 700 approach, the small radius leaves limited space around the object, making two-finger grasps more
 701 suitable; from an oblique side approach, the object height provides enough accessible surface for
 702 three-finger contacts, so the *Single-Three* prior favors oblique wrist poses.

703 **From Human Prior to Robot Grasps.** Fig. 21 visualizes how the learned human prior guides robot
 704 grasp synthesis. The sampled human-prior wrist poses provide coarse object-conditioned guidance,
 705 which is transferred to robot wrist initializations and then refined by force-closure-aware optimization.
 706 The examples cover both single-hand and bimanual grasps across different object scales.

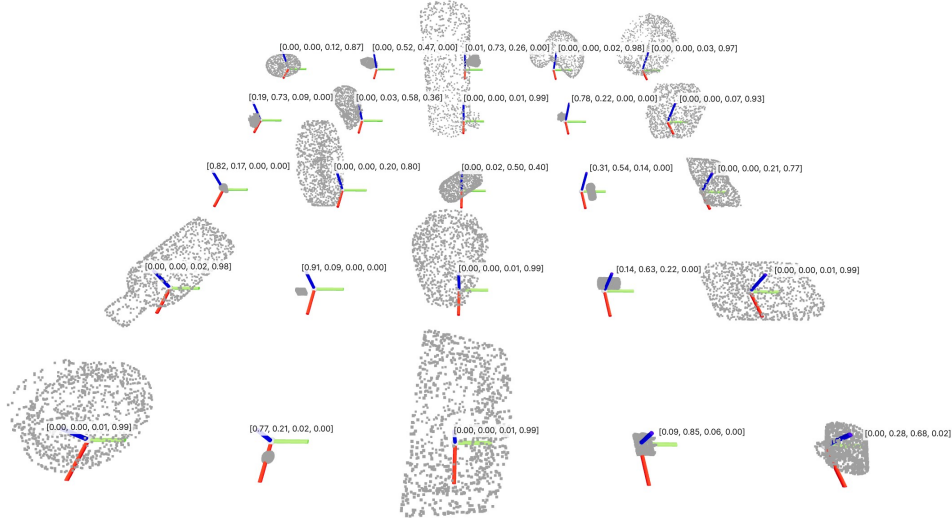


Figure 19: **Predicted contact-mode distributions from the human prior.** Each example shows an object point cloud together with the predicted probability vector over contact modes, ordered as [Single-Two, Single-Three, Single-Full, Both-Full]. The coordinate axes have a length of 0.1 m.

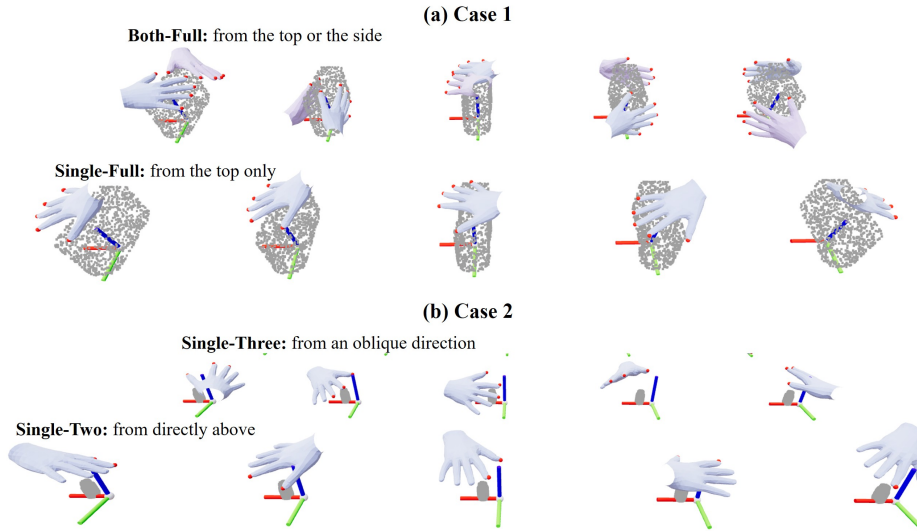


Figure 20: **Object-conditioned wrist-pose samples for different contact modes.** The examples show that the learned pose prior predicts wrist-pose distributions conditioned jointly on object geometry and contact mode. For a given object geometry, different contact modes favor different approach regions: in Case 1, *Both-Full* admits both top and side bimanual approaches, while *Single-Full* mainly concentrates on top-down approaches; in Case 2, *Single-Three* favors oblique approaches, whereas *Single-Two* favors more direct top-down approaches.

707 **Real-World Grasping Snapshots.** We provide real-world grasping snapshots in Fig. 22 and Fig. 23.
 708 The objects cover a wide range of scales and are approximately ordered from smaller to larger
 709 instances. Depending on the object geometry, each object may support one or two contact modes. For
 710 each feasible contact mode, we execute three distinct grasps to illustrate the diversity of generated
 711 bimanual grasp behaviors.

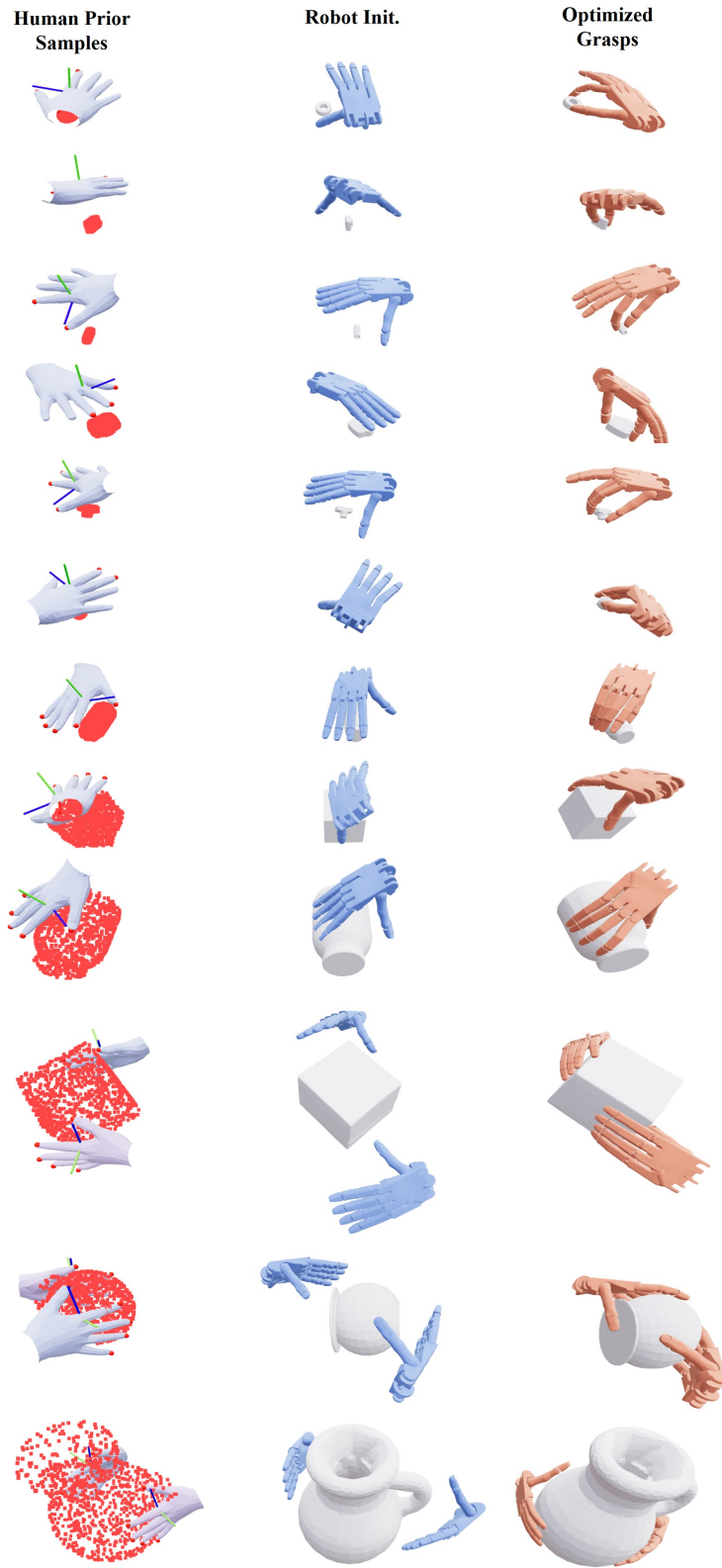


Figure 21: **From human prior samples to optimized robot grasps.** Left: contact-mode and wrist-pose samples predicted by the human prior for each object. Middle: robot wrist initializations obtained by transferring the prior samples to the robot hand. Right: final robot grasps after force-closure-aware optimization.



Figure 22: **Snapshots of real-world grasp on objects #1–#12.** Objects are approximately ordered by increasing scale from top to bottom. Some objects allow only one feasible contact mode, whereas others allow two. For each feasible contact mode, we execute three distinct grasps.

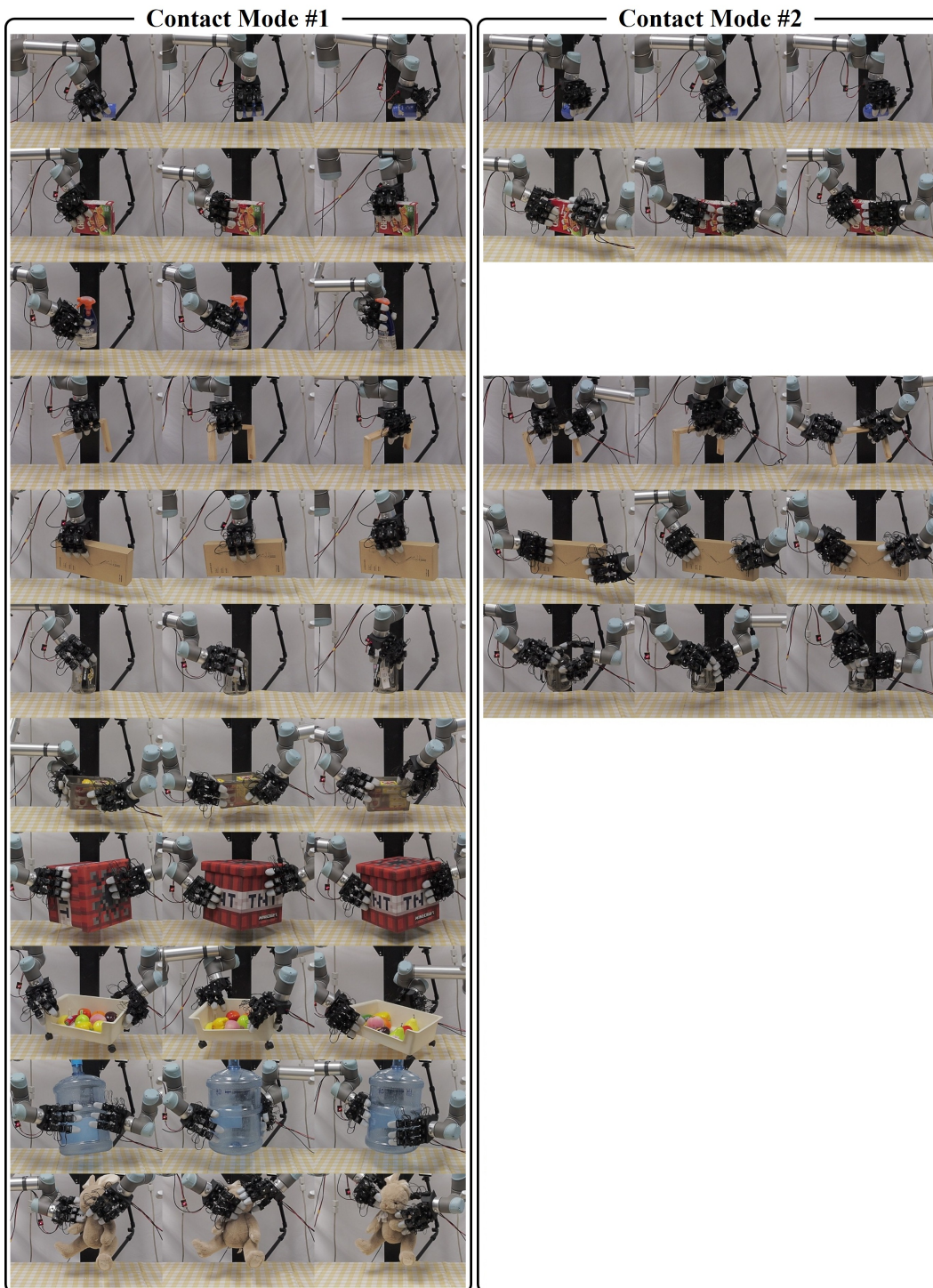


Figure 23: Snapshots of real-world grasp on objects #13–#23.

712 D Failure Analysis

713 D.1 Failure Analysis of Human-Prior-Guided Grasp Synthesis

714 Fig. 24 summarizes typical failure modes observed during human-prior-guided grasp synthesis.
715 Although the learned prior provides useful global guidance, it remains a coarse prior. First, the
716 human prior and the optimization objective do not explicitly model functional grasp awareness, so
717 the pipeline may produce grasps that are geometrically plausible but functionally weak. Second,
718 the contact-mode prior can sometimes predict *Single-Full* for objects that are slightly beyond the
719 capability of single-hand grasping. Third, for relatively small objects, the *Both-Full* prior can produce
720 partially overlapping hands. Finally, for extremely large objects, the *Both-Full* prior can produce
721 noticeable hand-object penetration. As described in Sec. A.3, the bimanual initialization offset helps
722 reduce such problematic initial states. These hand-hand overlap and hand-object penetration cases
723 can be further mitigated by the feasibility-constrained physical optimization stage. This highlights
724 the benefit of our human-prior-plus-robot-optimization design: the learned human prior provides
725 diverse and global grasp hypotheses, while the robot optimization stage explicitly enforces physical
726 plausibility and adapts the grasps to the target robot embodiment. These failure modes reveal
727 limitations of the current prior and highlight the need for larger and more comprehensive human
728 grasp datasets, together with stronger functional reasoning for grasp synthesis.

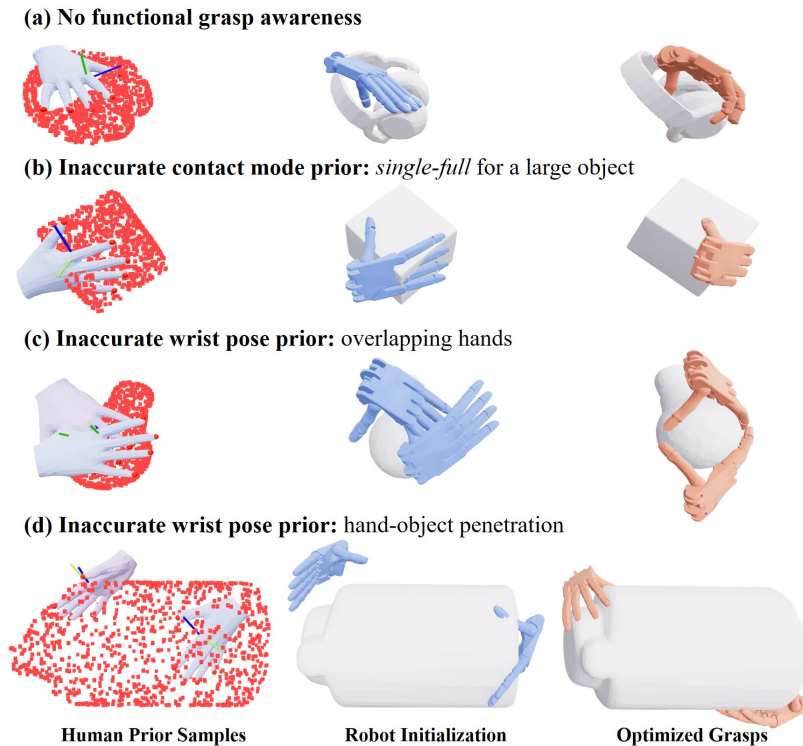


Figure 24: **Failure analysis of human-prior-guided grasp synthesis.** Each row shows a failure mode from human prior samples to robot initialization and optimized grasps. (a) The human prior and the optimization objective lack functional grasp awareness, so the pipeline may produce geometrically plausible but functionally weak grasps. (b) The contact-mode prior can sometimes predict *Single-Full* for objects that are slightly beyond the capability of single-hand grasping. (c) For relatively small objects, the *Both-Full* prior can sometimes produce partially overlapping hands. (d) For extremely large objects, the *Both-Full* prior can sometimes produce clear hand-object penetration. The hand-hand overlap and hand-object penetration cases can often be mitigated during optimization.

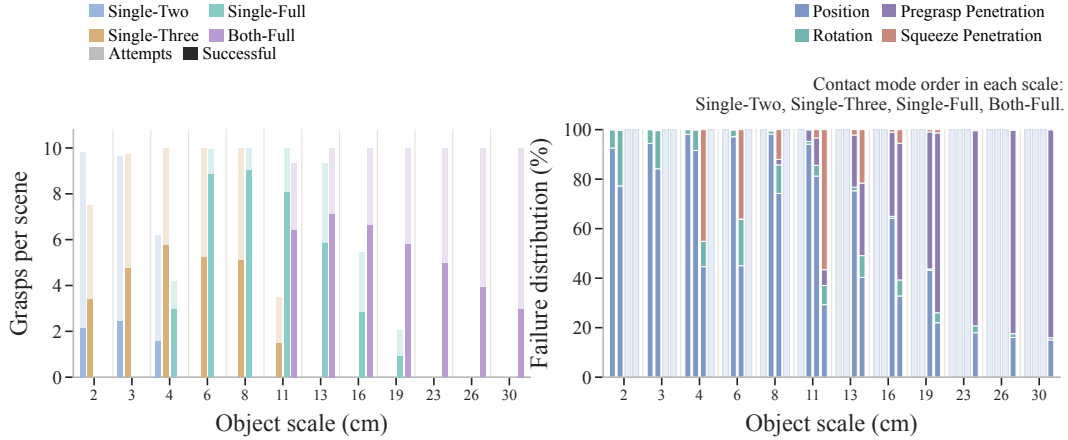


Figure 25: **Failure distribution of the distilled grasp generator on the Shadow Hand.** Left: number of generated and successful grasps per scene across object scales and contact modes. Right: distribution of failure causes for each contact mode at each object scale, ordered as *Single-Two*, *Single-Three*, *Single-Full*, and *Both-Full*. Failures are categorized as position error (unsuccessful lifting), rotation error (excessive object rotation), pregrasp penetration, and squeeze-stage self-penetration.

729 D.2 Failure Analysis of Distilled Robot Grasp Generators.

730 Fig. 25 further analyzes generator failures after distillation by reporting both the simulation success
 731 rates of generated grasps and the distribution of failure causes across object scales and contact
 732 modes. The results reveal several scale- and mode-dependent failure patterns. First, most failures
 733 are caused by position error, where the generated grasp does not successfully lift the object. This
 734 indicates a remaining limitation in the accuracy of grasp-pose generation. For *Single-Full* grasps,
 735 a noticeable portion of failures comes from squeeze-stage self-penetration, which is not explicitly
 736 considered in the current squeeze-pose synthesis process [2]. This issue is particularly pronounced
 737 on the Shadow Hand, whose relatively small inter-finger spacing makes self-penetration more likely
 738 during finger closure. We plan to address this limitation in future work. For large objects, especially
 739 under the *Both-Full* mode, many failures already occur at the pregrasp stage due to hand-object
 740 penetration, suggesting that the generated global poses are not sufficiently accurate for large-scale
 741 objects. These observations indicate that, although HUGS-synthesized data can effectively supervise
 742 online generators, designing robust robot grasp generators across scales and contact modes remains
 743 an important direction for future research.

744 D.3 Failure Analysis of Real-World Grasping

745 In addition to the generative-network limitations analyzed in the simulation evaluation, real-world
 746 failures are affected by several deployment-specific factors, as illustrated in Fig. 26. First, low-quality
 747 grasp predictions can produce loose grasps that fail to securely lift the object. Second, calibration,
 748 segmentation, and depth reconstruction errors can corrupt the reconstructed object point cloud,
 749 causing the generated grasp to miss the target object, especially for small objects. Third, geometrically
 750 irregular objects and challenging cross-scale cases may lie outside the training distribution, leading to
 751 degraded grasp generation. These observations highlight the need for advances at both the model
 752 and system levels. On the model side, future work should improve the accuracy and generalization
 753 of grasp generators to more diverse object distributions. On the system side, robust real-world
 754 deployment will require end-to-end closed-loop execution, tactile adaptation, and collision-aware
 755 motion planning.

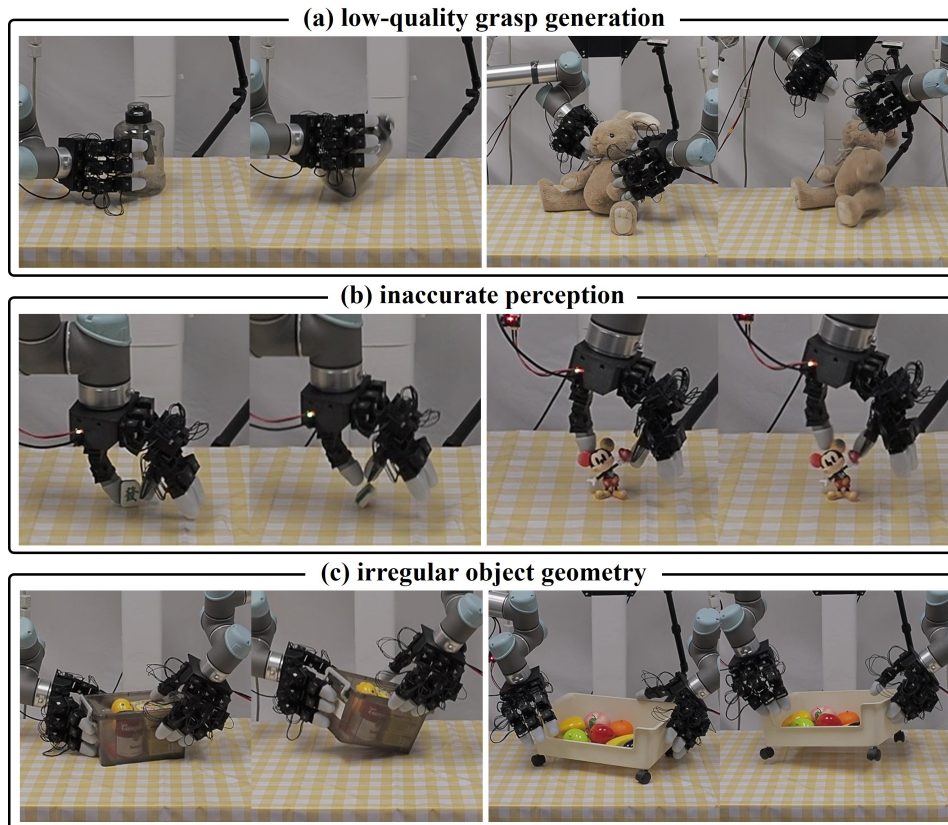


Figure 26: **Typical failure modes of real-world grasping.** (a) Low-quality grasp generation may produce loose or unstable grasps, suggesting that the current grasp generation network can be further improved. (b) Inaccurate calibration, segmentation, and depth sensing reduce the reliability of reconstructed point clouds, resulting in grasps that miss the object. This issue is particularly pronounced for small objects. (c) Irregularly shaped objects often suffer from poor point-cloud reconstruction and may lie outside the distribution of the training object set, leading to degraded grasp generation. This suggests that the diversity and coverage of the synthetic object dataset can be further improved to enhance generalization.