

# DexGraspNet 2.0: Learning Generative Dexterous Grasping in Large-scale Synthetic Cluttered Scenes: Supplementary materials

Anonymous Author(s)

Affiliation

Address

email

1 This document provides supplementary details, additional experiments, and enhanced visualizations  
2 to complement the main paper. Sec. 1 outlines the detailed experimental settings discussed in the  
3 main paper. Sec. 2 presents additional experiments conducted to extend the findings. Sec. 3 offers  
4 further statistical insights into our DexGraspNet 2.0 benchmark. Sec. 4 details the methodology  
5 used for generating grasp labels. Sec. 5 elaborates on the technical aspects involved in constructing  
6 and training our model for dexterous grasping in cluttered scenes. Sec. 6 highlights specific imple-  
7 mentation details related to applying our model with parallel grippers. Sec. 7 provides additional  
8 visualizations showcasing our dataset and model.

## 9 1 Experiment Details

10 We provide additional details on the experiment settings due to space constraints in the main paper.  
11 Sec. 1.1 delineates how we evaluate a grasp in a simulator and enumerates some of the physics  
12 parameters involved. Sec. 1.2 elaborates on the three ablation groups in detail. Sec. 1.3 outlines the  
13 three baseline methods benchmarked in the main paper.

### 14 1.1 Evaluation Metric

15 We evaluate various grasping models by measuring their simulation success rates in the Isaac Gym  
16 simulator. For each test scene, a model is expected to take a single-view depth point cloud as input  
17 and output one grasp pose  $G_p$ . If capable of generating multiple grasps, the model must select the  
18 best proposal, as required in the main paper. Following this, the evaluator determines whether  $G_p$   
19 constitutes a successful grasp. Specifically, a predefined rule is applied to calculate a pregrasp pose,  
20 squeeze pose, and lift pose, thereby establishing a complete action trajectory  $T$ . Subsequently,  $T$  is  
21 executed within the simulator, and success is determined by its ability to lift an object off the table  
22 without any initial intersection with the table or surrounding objects. Consistency is ensured across  
23 all experiments by maintaining the same trajectory generation rule and physics parameters. Some  
24 of the important physics parameters are listed in Tab. 1.

Parameter	Value	Parameter	Value
friction coeff	0.2	object mass	0.1 kg
joint stiffness	800	joint damping	20

Table 1: Physics Parameters

### 25 1.2 Ablation Details

26 We explain the settings of the three ablation groups from our main paper in detail.

27 **Local Feature.** Our grasping method aims to achieve higher generalization efficiency by condi-  
28 tioning on local features. We investigate this design by training a diffusion model that predicts the

distribution of all valid grasps conditioned on the scene’s global feature. This ablated version has three major differences compared to our original model: (1) it discards the UNet decoder and retains only the encoder; (2) during training, each grasp label corresponds to the global feature vector of the scene point cloud (output by the encoder) instead of the local feature vector of the grasp’s corresponding point (one of the point-wise vectors output by the decoder). (3) during inference, the model does not predict graspness or propose seed points, but only encodes the scene point cloud and directly generates grasp poses conditioned on its global feature vector.

**Decomposed Pose Modeling.** Our grasp generation module models the conditional distribution  $p(T, R, \theta|f_s)$  in a decomposed manner: a conditional generative model predicts the conditional distribution  $p(T, R|f_s)$ , followed by a deterministic model predicting  $\theta$  from  $f_s$  and  $(T, R)$ . Surprisingly, the above design slightly outperforms a seemingly more elegant approach: using a single conditional generative model to fit the joint distribution  $p(T, R, \theta|f_s)$  without decomposing the wrist pose  $(T, R)$  from the joint angles  $\theta$ . We postulate that this phenomenon results from the distribution of the training data, rather than an inability to properly tune the second approach. Specifically, our training dataset primarily consists of power grasps that utilize all fingers, resulting in a single-mode distribution of  $\theta$  conditioned on  $(T, R)$  and  $f_s$ . Consequently, the deterministic model regressing  $\theta$  is not confused by this data distribution; instead, it potentially becomes more robust to outliers. Essentially, the outcomes of this ablation group are highly specific to our task and training data. If we incorporate additional grasping modes into our dataset, such as precision grasps and functional grasps, it would violate our assumption of a single-mode distribution of  $\theta$  conditioned on  $(T, R)$  and  $f_s$ . In such a scenario, jointly modeling  $p(T, R, \theta|f_s)$  with a single conditional generative model might outperform our current design.

**Randomly-Packed Training Scenes.** In addition to ablating our network designs, we also conduct one experiment to ablate our dataset in the main paper. Our training set comprises 100 densely-packed scenes (with 8 to 11 objects) and 7500 randomly-packed scenes (with 1 to 10 objects). All dense scenes are sourced from [1]. However, we observed that training solely on these dense scenes resulted in the inability to generate valid grasp poses when the table is nearly clear. Therefore, we incorporated the randomly-packed scenes to ensure performance across all density levels.

### 1.3 Baseline Details

We outline the three baselines compared in the main paper and detail how we adapted two of them from their original setting of single-object grasping to our cluttered scenarios.

**HGC-Net [2].** HGC-Net is a two-stage method for grasping in cluttered scenes. Initially, a segmentation model divides the scene point cloud into graspable points and ungraspable points. Following this, a deterministic model predicts a grasp pose near each graspable point. Given that this method already focuses on cluttered scenes, minimal modifications were required. The only change made was switching their end effector from the HIT-DLR II hand to the LEAP hand.

**ISAGrasp [3].** ISAGrasp is a regressive method designed for grasping single objects. It employs a PointNet++ encoder [4] to encode the object point cloud into a global feature vector. Subsequently, an MLP is utilized to predict the wrist translation, wrist quaternions, and joint angles. We extensively modified this method to adapt it for cluttered scenes: (1) We replaced their PointNet++ encoder with a ResUNet14 encoder-decoder and incorporated a seed point proposal module based on point-wise graspness prediction, similar to our method. (2) During inference, this modified model predicts the grasp parameters from the local feature vector of the proposed seed point, instead of the global feature vector obtained from their original point cloud encoder. (3) During training, each grasp label is associated with its corresponding point rather than its target object. We designate the modified model as ISAGrasp<sup>†</sup>. It is worth noting that this adaptation already rectifies a major suboptimal aspect of their original baseline by integrating one of our key designs: replacing global conditioning with local conditioning. Consequently, the adapted method differs from our model solely in the use of a regressive model to predict the wrist pose, whereas we employ a conditional generative model.

	Method	GraspNet-1Billion			ShapeNet		
		Dense	Random	Loose	Dense	Random	Loose
Ablation	Euler Angle	87.6	82.0	73.0	78.0	76.4	<b>75.2</b>
	Axis Angle	86.4	81.7	70.5	79.0	76.4	74.1
	Quaternion	87.9	81.5	72.0	78.6	77.0	72.9
	6D	88.2	81.5	71.9	80.2	79.0	73.0
	Ours	<b>90.6</b>	<b>83.7</b>	<b>73.2</b>	<b>81.0</b>	<b>85.4</b>	74.2

Table 2: **Ablation studies for representations of rotation.** **Euler Angle** represents rotation as 3D Euler angle; **Axis Angle** represents rotation in 3D as the angle of rotation multiplies the rotation axis; **Quaternion** represents rotation as 4D quaternion; **6D** represents rotation with the first two rows of the rotation matrix. **Ours** represents the rotation as the rotation matrix.

	Method	GraspNet-1Billion			ShapeNet		
		Dense	Random	Loose	Dense	Random	Loose
Ablation	Graspness	81.8	76.6	68.0	73.7	71.3	64.4
	Log Probability	78.1	78.4	75.1	72.4	71.6	<b>74.6</b>
	Random	65.1	62.0	57.2	61.7	58.9	56.4
	Ours	<b>90.6</b>	<b>83.7</b>	<b>73.2</b>	<b>81.0</b>	<b>85.4</b>	74.2

Table 3: **Ablation studies for sampling strategy.** **Graspness** ranks samples by graspness score only; **Log Probability** ranks samples by log probability only; **Random** randomly draws from sampled poses; **Ours** ranks samples by combination of graspness scores and log probabilities.

78 **GraspTTA [5].** GraspTTA utilizes a CVAE for grasping single objects. It leverages PointNet [6]  
79 to encode the object point cloud into a global feature vector, which serves as conditioning for the  
80 CVAE to predict the distribution of the wrist translation, wrist axis angles, and joint angles. We  
81 adapt it for cluttered scenes using the same approach as ISAGrasp<sup>†</sup>, and denote the adapted version  
82 as GraspTTA<sup>†</sup>. Furthermore, we discard the test-time optimization of the original method because  
83 it relies on the full point cloud, which is an invalid assumption in our task settings.

## 84 2 Additional Experiments

### 85 2.1 Ablate Rotation Representation

86 Our method employs the rotation matrix to represent wrist rotation and applies SVD [7] to orthog-  
87 onalize network predictions. We compared this design against several alternatives: **Euler Angle**  
88 (representing rotation as 3D Euler angles), **Axis Angle** (rotation represented by the angle of rota-  
89 tion multiplied by the rotation axis), **Quaternion** (represented as a 4D quaternion), and **6D** (using  
90 the first two rows of the rotation matrix). The results in Tab. 2 demonstrate that our choice out-  
91 performs all other methods across the evaluated task. As discussed in [7], rotation representations  
92 in Euclidean space with fewer than five dimensions, such as Euler angles, axis-angle, and quater-  
93 nions, are inherently discontinuous. Although the 6D representation circumvents this issue, it is  
94 coordinate-dependent. Introducing small noises in different directions to the rotation in a 6D rep-  
95 resentation results in changes of varying magnitudes. In contrast, our 9D representation is both  
96 continuous and coordinate-independent, thereby outperforming other rotation representations.

### 97 2.2 Ablate Ranking Strategy

98 During inference, we rank all predicted samples to identify the best one using a linear combination  
99 of the graspness scores of the seed points and the estimated log probabilities of the wrist poses. We  
100 ablate this ranking strategy by removing the graspness score, the log probability, or both. Tab. 3  
101 presents the results. Our method (**Ours**), which ranks samples based on a combination of grasp-  
102 ness scores and log probabilities, consistently outperforms the other strategies. Ranking solely by  
103 graspness scores (**Graspness**) or log probabilities (**Log Probability**) yields moderate performances,  
104 while selecting samples randomly (**Random**) results in the lowest success rates. These findings un-  
105 derscore the efficacy of our proposed ranking strategy in identifying optimal grasp poses.

106 Interesting to note, despite the theoretical challenges in defining a probability density function  
 107  $p(T, R|f_s)$  on a 6-dimensional data manifold embedded within a higher-dimensional parameter  
 108 space (12D), experiments demonstrate that our estimated log probabilities consistently enhance the  
 109 performance of our ranking strategy. Nevertheless, we acknowledge this theoretical inelegance and  
 110 defer the solution to future studies, such as exploring the use of normalizing flows on  $SE(3)$  or  
 111 employing manifold diffusion methods.

### 112 2.3 Scaling the Dataset for Grippers

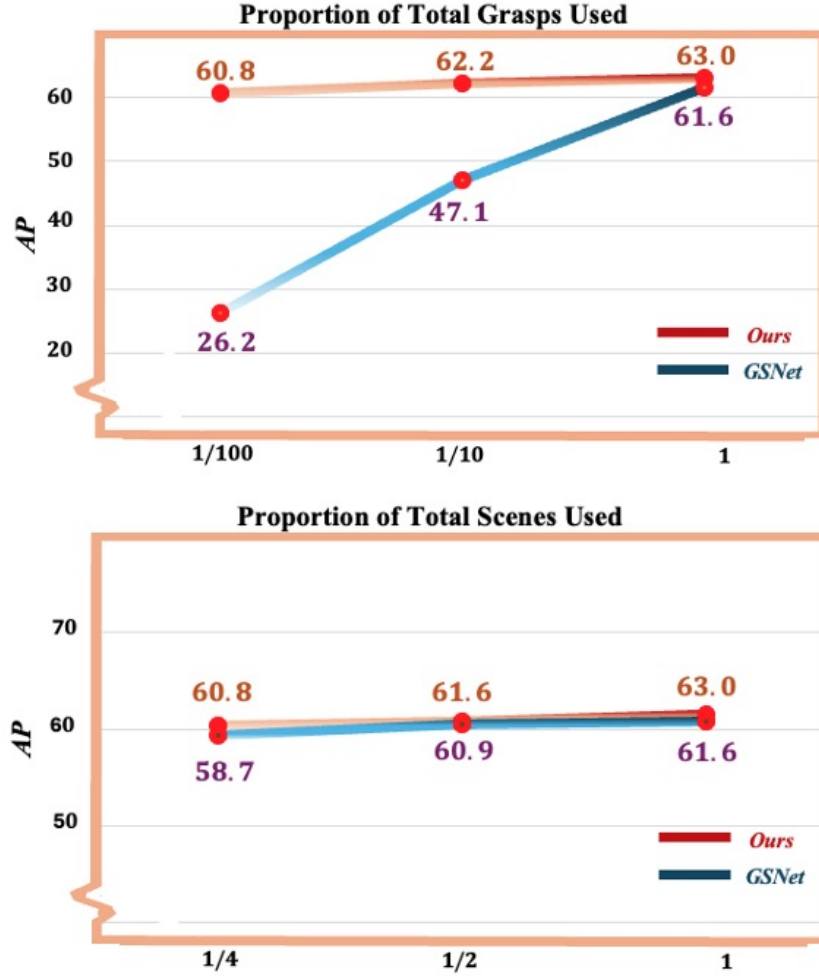


Figure 1: **AP metric** evaluated on models trained with downscaled dataset. **Top**: downsampling the number of grasp labels in each scene. **Bottom**: downsampling number of scenes trained on

Fraction of Grasps	Success Rate
1/100(42k)	81.3
1(4.2M)	92.4

Table 4: **Success Rate of real-world experiment on Ours model trained over downscaled dataset.** We train Ours model with random 1/100 fraction of grasp labels and the entire grasp pose dataset, amounting 42k and 4.2M labels, respectively.

Depth Restoration	Diffuse	Trans	Hybrid
With	94.1	80.0	90.7
Without	94.1	50.0	86.4

Table 5: Real-world cluttered scene dexterous grasping with/without depth restoration. **Diffuse** includes only diffuse objects, **Trans** comprises only transparent or specular objects, and **Hybrid** includes scenes used in the main paper, consisting of a mixture of diffuse, transparent, and specular objects for comparison.

End Effector	Normal	Large
Parallel Gripper	92.4	0.0
Dexterous Hand	81.5	100.0

Table 6: Comparison of real-world grasping performance using a parallel gripper or a dexterous hand across different scene types. The five **Normal** scenes consist of typical cluttered environments, while the **Large** scene includes 4 large objects.

113 We scale down the training data of parallel gripper by (1) reducing the number of grasps in each  
 114 scene, and (2) decreasing the number of training scenes. We evaluate the AP metric in simulation  
 115 for each setting and success rate in real world.

116 As shown in Fig. 1, although under the full-data setting our generative model only slightly outper-  
 117 forms GSNet by +1.4 AP, the AP metric of GSNet drops by a significant amount of 35.4 as we  
 118 downscale the number of grasps by 100, whereas our generative pipeline drops by only 2.2. This  
 119 suggests that our generative pipeline is significantly more sample-efficient than GSNet. Both meth-  
 120 ods are robust to downscaling of number of training scenes at the scope of our experiment, with only  
 121 slightly dropped AP.

122 The resulting statistics in terms of AP is much to our surprise, as being trained with 1/100 total grasp  
 123 labels, namely only 42k grasp labels, our generative model seems to still retain strong performance.  
 124 In order to validate this counter-intuitive result, we carry out real-robot experiments with Ours mod-  
 125 els trained with downsampled number of grasps and report success rate in Tab. 4. With 42k training  
 126 labels, our generative model achieve 81.5% success rate in real-world cluttered scenes as shown in  
 127 Fig. 2, which is affirmative to the AP statistics.

128 In summary, the experiments in this section give strong evidence that the distribution of valid grasp  
 129 poses does exist and the amount of data required to simulate at least a valid support of such a  
 130 distribution may prove to be much smaller than previously been conjectured.

## 131 2.4 Using Raw Depth in the Real World

132 In our real-world experiments, we integrated depth restoration techniques [8] to facilitate grasping  
 133 transparent and specular objects amidst cluttered scenes. Here, we conduct additional experiments  
 134 to demonstrate that our method do not rely on depth restoration when grasping diffuse objects.  
 135 We constructed four additional cluttered scenes in the real world: two scenes (**Diffuse**, as shown in  
 136 Fig. 2) consisting solely of diffuse objects and two scenes (**Trans**, as shown in Fig. 3) containing only  
 137 transparent and specular objects. The original five test scenes from the main paper, which include  
 138 a mixture of objects, are denoted as **Hybrid**. We then evaluated our model on all test groups both  
 139 with and without the application of depth restoration techniques. The results in Tab. 5 demonstrate  
 140 two key findings: firstly, our model’s effectiveness in real-world grasping is independent of depth  
 141 restoration for **Diffuse** scenes; secondly, our model exhibits enhanced robustness to object texture,  
 142 particularly transparent and specular surfaces, when depth restoration is applied.

## 143 2.5 Discussion on Dexterous Hands vs Parallel Grippers

144 While grasping systems utilizing parallel grippers have already achieved impressive robustness in  
 145 the real world [9, 10], we advocate that dexterous hands can further enhance performance. In ad-



dition to the 5 test scenes (**Normal**, as shown in Fig. 4) demonstrated in the main paper, we also construct an additional scene (**Large**) consisting of 4 large objects, as shown in main paper. Real-world experiment results in Tab. 6 indicate that the dexterous hand can grasp each object in this scene, whereas the parallel gripper cannot grasp any object. This is because the dexterous hand possesses strong envelopment capabilities, allowing it to grasp larger objects effectively.

### 151 3 Benchmark Specifications

This Section presents further details about the DexGraspNet 2.0 benchmark proposed by this work. Sec. 3.1 provides statistics of the DexGraspNet 2.0 benchmark, including both the **Training Set** that contains ground truth grasp pose annotations and the **Test Set** with no ground truth provided. Sec. 3.2 identifies the objects used to generate our benchmark. Sec. 3.3 presents the pipeline used to generate training scenes with selected objects. Sec. 3.4 elaborates the protocol of generating test scenes and how we divide them into different splits.



### 3.1 Benchmark Statistics

Splits	number of objects	number of scenes
Training	60(GraspNet1B)	100(seminal)+7500(augmented)
Test	88(GraspNet1B) + 1231(ShapeNet)	670
Total	88(GraspNet1B) + 1231(ShapeNet)	8270

Table 7: Statistics of the DexGraspNet 2.0 Benchmark

Tab.7 illustrates the overall statistics. The entire benchmark encompasses two components: a **Training Set** used to train our models and a **Test Set** to evaluate dexterous grasping pose generation models on. Note that ground truth grasp pose annotations are only provided for training set. In total, the benchmark contains 8270 scenes, 1319 objects and 426.6M grasp pose annotations.

**Training Set** contains 7600 scenes and 60 objects in total. all training objects are from the GraspNet-1Billion [1] dataset

**Test Set** contains 670 scenes and 1319 objects in total. the 88 objects from the GraspNet-1Billion [1] dataset are used to compose 450 of the test scenes, and 1231 objects picked from ShapeNet [11] are used to compose the remaining 220 test scenes

### 3.2 Object Selection

The 60 objects in Training Set are those appeared in GraspNet-1Billion [1] scenes 0000-0099. The Test Set contains 1319 objects, 88 of them are all the objects in GraspNet-1Billion [1], and the remaining 1231 objects are picked from ShapeNetSem [11].

### 3.3 Training Scenes Specification

In the 7600 training scenes, 100 are called **seminal scenes**, which corresponds to the Scenes 0000-0099 in the GraspNet-1Billion [1] dataset composed and rendered using their official meshes and annotations. We augment each seminal scenes 75 times by randomly deleting objects in the scene. In each augmented scene, the number of objects deleted is uniformly sampled from  $[1, k-1]$ , where  $k$  is the number of objects in the original scene. In total, we generate 7500 augmented training scenes with 100 seminal scenes, totalling 7600 scenes in the entire training set.

### 3.4 Test Set Scenes Specification

As shown in Tab. 1 of the main paper, the Test Set is divided into 6 splits. In the following, we specify each of these splits.

**GraspNet-1Billion Dense** composes of 90 scenes that correspond to the Scenes 0100-0189 in the GraspNet-1Billion [1] dataset. Each scene contains 8-11 objects.

**GraspNet-1Billion Random** composes of 180 scenes. This split is generated by augmenting each GraspNet-1Billion Dense split scenes twice with the process as described in Sec.3.3

**GraspNet-1Billion Loose** composes of 180 scenes by augmenting each GraspNet-1Billion Dense split scenes twice with the process as described in Sec.3.3, with only 1-2 random objects remaining in the scene.

The three ShapeNet splits are generated by dropping objects on a 30cm×50cm tabletop. In specific, we follow the scene generation process of DREDS [12] with the material randomization function disabled. We run the scene generation process in PyBullet [13] and filter physically stable ones in IsaacGym [14]. The Dense/Random/Loose splits are divided according to the number of objects appearing in each scenes.

**ShapeNet Dense** composes of 100 scenes, each containing 8-11 objects

**ShapeNet Random** composes of 90 scenes, each containing 5-9 objects

196 **ShapeNet Loose** composes of 30 scenes, each containing 1-2 objects

## 197 **4 Grasp Label Generation**

198 This section elaborates our pipeline for generating dexterous grasping poses on single objects. First,  
199 we define initial hand poses by retargeting GraspNet-1Billion [1] annotations to dexterous hand.  
200 Then we run physics-based optimization to generate stable grasps. To maximally diversify the pro-  
201 duced data, we adopt two different methods, [15] which targets Grasp Wrench Space (GWS) opti-  
202 mality, and [16] which targets force-closure, as optimization algorithms, each generating half of the  
203 dataset. Lastly, we filter stable and collision-free grasps via simulation in the IsaacGym [14] simu-  
204 lator. As shown in Fig. 5, in total we generate 44.9M stable grasp poses for 88 objects from 280M  
205 initial poses. Even in the face of our very strict friction coefficient  $\mu=0.2$ , our method still maintains  
206 overall success rate of 16.07%. In the following subsections we detail each of these components.

### 207 **4.1 Hand Pose Initialization**

208 As discovered in [16], the success rate of dexterous grasp generation is very sensitive to initial hand  
209 pose. Moreover, we aim to cover valid grasp modes for each object as comprehensively as possible.  
210 Therefore, we initialize dexterous hand poses by retargeting the exhaustive GraspNet-1Billion [1]  
211 gripper annotations.

212 In specific, we filter points where stable gripper grasp poses are annotated in [1] as grasp points. As  
213 shown in Fig. 7, for each grasp point, we align the +y axis (pointing forward out of the palm) of  
214 dexterous hand with the +x axis of gripper pose annotation, retreat the center of palm a fixed distance  
215 from grasp point in the approaching direction, initialize hand joint qpos with a set of predefined values  
216 and exhaustively apply transformations corresponding to 256 approaching directions, 4 depths and  
217 12 in-plane angles as defined in [1].

### 218 **4.2 Grasp Pose Optimization**

#### 219 **4.2.1 GWS-based optimization (adapted version of [15])**

220 We reimplement [15] on the CuRobo [17] framework for better computation parallelism. We set the  
221 target Task Wrench Space (TWS) as a unit sphere in 6D wrench space such that the task objective is  
222 identical to forming a force-closure grasp, and run 600 iterations with naive gradient descent.

#### 223 **4.2.2 force-closure-based optimization (adapted version of [16])**

224 We adopt [16] with modification in its definition of force-closure energy, and reimplement the mod-  
225 ified algorithm on the CuRobo [17] framework as well.

226 We observe that the force-closure energy used in [16] assumes unit contact force is applied to each  
227 contact point, whereas human naturally adjust contact forces applied to different contact points in  
228 order to maintain a firm grasp. The above assumption limits the objective of optimization in [16]  
229 onto a submanifold of the space of all valid grasp poses, hurting the quality and diversity of generated  
230 data. Following the notations in [16], we relax the unit-contact force assumption by reformulating  
231 the force closure energy as the following bilevel form:

- 232 • At each timestep, given the current hand pose, we solve the optimal contact forces applied to  
233 current contact points such that the total wrench imposed on the object is minimized. We formulate  
234 this intuition into the following linear program:

$$\begin{aligned} P_t &= \min_{\lambda_t} \|G(\lambda_t \odot c)\|_2 \\ s.t. & \max_i (\lambda_t)_i = 1 \\ & (\lambda_t)_i \geq 0, i = 1, 2, \dots, n \end{aligned}$$



Where  $P_t$  has the physical meaning as the total wrench applied to the object when the combination of contact force magnitude,  $\lambda_t$ , is applied to the contact points.  $\odot$  means element-wise product. Note this linear program admits closed-form solution therefore imposes neglectable computation burden.

• Across timesteps, we optimize the differentiable force-closure metric in awareness of the plausibility of the current hand pose:

$$E_{FC} = \begin{cases} \|G(\lambda_t \odot c)\|_2, & \text{if } P_t < \tau_{FC}, \min_i(\lambda_t)_i \geq \tau_\lambda, \text{ and } B = 1 \\ \|Gc\|_2, & \text{otherwise} \end{cases}$$

Where  $\tau_{FC}, \tau_\lambda$  are predefined thresholds, and  $B$  is a binary random variable with  $P(B = 1) = 0.9$ .

If the current hand pose is already capable of forming a force-closure grasp on the object, mathematically defined as  $P_t < \tau_{FC}$  (total wrench acceptably small) and  $\min_i(\lambda_t)_i \geq \tau_\lambda$  (a minimum contact force is applied to each contact point), then we decide the current pose is good enough in terms of force-closure property. In this case, we scale the force closure energy to prevent overoptimization. In effect, the force closure energy now works as a regularization term. Otherwise, if the current hand pose is not stable enough, we keep searching for more stable poses by optimizing the force closure metric with original energy term. In addition, even for the former case, we stochastically use the original energy term with probability 0.1 to encourage forming more robust grasp poses.

Note in the above formulation, the global minimum set of hand poses for  $E_{FC}$  are the poses for which there exists a non-trivial contact force combination such that the total wrench executed to the object is zero. This global minimum set exactly corresponds to the original definition of force closure in [18].

### 4.3 Filtering Stable and Collision-Free grasps

We perform grasp filtering in the IsaacGym simulator. First, we check for each grasp pose if the penetration between hand mesh and object mesh is below 2 mm. For all collision-free grasps, we execute the grasp with a predefined heuristic and simulate for 60 timesteps at 60Hz. The grasp pose is validated as stable if it can deny gravity in all 6 axis-aligned directions. The friction coefficient  $\mu$  for both hand and objects are set to 0.2, making the filtering process very strict.

Fig. 5 shows the **Valid Rate** for each object, which is defined as the portion of generated grasps that are both collision-free and stable. The overall success rate is 16.07%, as we generate in total 44.9M valid grasp poses out of 280M grasp pose initializations. The method-specific valid rate for [15] and [16] are 7.91% and 24.19% respectively.

## 5 Implementation Details for Dexterous Hands

In this section, we elaborate on the data organization (Sec. 5.1) and model architecture (Sec. 5.2) of our method for dexterous grasping.

### 5.1 Data

**Data Reblancing** In each training scene, the numbers of grasp labels on graspable objects may be uneven. Randomly sampling grasp labels uniformly across all valid ones in each scene could slow down the learning of grasping objects that have fewer labels. To address this, we implement a two-stage sampling approach to rebalance the training process: first, we randomly sample a graspable object, and then we randomly sample one of its labels.

**Data Augmentation.** We implement data augmentation by rotating the scene point cloud and grasp labels around the camera axis with a random angle uniformly sampled from the interval  $[0, 2\pi)$ . No further augmentations are needed.

**Ground-truth Graspness Definition.** For each training scene, we define a graspness score for the surface points of each object to represent its graspability. This score is determined by identifying a seed point and then assigning graspness to the nearby points. For an object  $o$  in this scene, we denote all valid grasp labels that target  $o$  as  $G_o = \{g_o^i\}$ , and the surface points of  $o$  as  $P_o = \{p_o^j\}$ . We then define a grasp cone with  $c$  being the apex, vector  $cm$  being the axis and an aperture of  $60^\circ$ , as shown in Fig. 8. Subsequently, we compute the projected distance of vector  $cp_o^j$  along  $cm$ , denoted as  $d$ , and the spanning angle  $\theta$  between  $cp_o^j$  and  $cm$ . Using these quantities, the value of  $f(g_o^i, p_o^j)$  is defined in Eq. 1. Numerically, this function is designed to attenuate exponentially with response to  $\theta$  and  $d$ , halving at  $10^\circ$  or 1.5 cm. Then the seed point is defined as the point with the largest  $f$  as shown in Eq. 2.

Finally, the seed point assigns graspness to nearby points with exponential decay and the graspness score of  $p_o^j$  is computed as the logarithm of the sum of all contributed graspness, as in Eq. 4. Empirically, this score reflects the number of valid grasp labels near  $p_o^j$ .

From another perspective, this correspondence implicitly defines a grasp distribution conditioned on a point within a scene. Although articulating this distribution in precise mathematical terms is difficult, we contend that it objectively exists. This distribution represents the target distribution that the grasp generation module approximates.

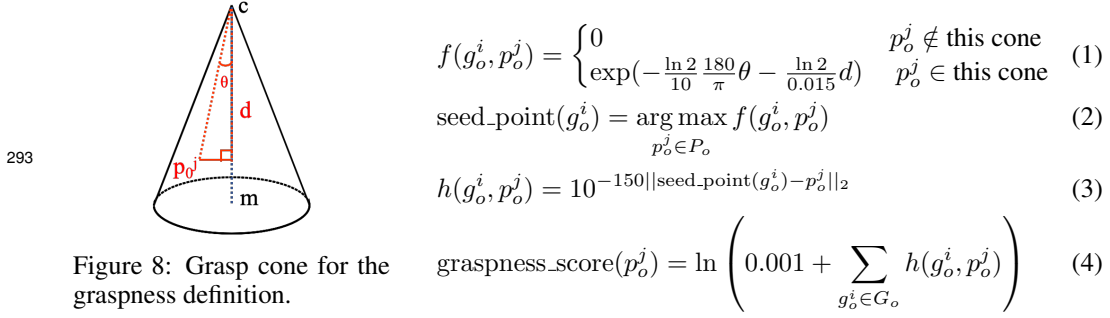


Figure 8: Grasp cone for the graspness definition.

## 5.2 Model

**Network Structure.** In the following paragraph, we elaborate on the network structures of our feature extractor, denoising model, graspness MLP, and joint MLP. First, our feature extractor employs the ResUNet14 architecture implemented with MinkowskiEngine [19] to derive point-wise feature vectors  $f_p \in \mathbb{R}^{512}$  from a scene point cloud  $P$ , which is quantized into sparse voxels. This network resembles the one utilized in GSNet [9]. Second, our denoising model  $v_\Theta(\hat{g}_E^t, f_s, t)$  is implemented as an MLP with layer sizes (524, 512, 256, 12) and Mish activations [20]. This model embeds  $t$  into  $\mathbb{R}^{512}$  using sinusoidal position embedding, adds this embedding with  $f_s$ , concatenates the resulting sum with  $\hat{g}_E^t$ , and feeds this concatenation into the MLP to predict the velocity. Third, our graspness MLP comprises a single-layer linear transformation, which maps  $f_p$  to three values. The first two are interpreted as binary classification logits indicating whether this point is an object point, while the third value represents the predicted graspness score  $GP_p$ . Fourth, our joint MLP is a 6-layered MLP with ReLU activations and residual block designs following [21].

**Detailed Diffusion Dynamics.** The forward and backward processes of the diffusion each consist of  $T_{\text{train}}$  and  $T_{\text{inference}}$  time steps, respectively, evenly distributed within the interval  $[0, 1]$ . Additionally, the number of time steps of the backward process is required to be a divisor of that of the forward process. We denote the interval between two neighboring time steps of the backward process as  $dt = 1/T_{\text{inference}}$ . The DDPM [22] scheduler is employed to schedule the forward process variances  $\beta_t$  for each time step  $t = i/T_{\text{train}}, i = 1, 2, \dots, T_{\text{train}}$ :

$$\beta_t = \beta_{\min} + \frac{i-1}{T_{\text{train}}-1} (\beta_{\max} - \beta_{\min}) \quad (5)$$

where  $\beta_{\min}, \beta_{\max}$  are hyper-parameters. Then we define  $\alpha_t = 1 - \beta_t$  and its cumulative product as  $\bar{\alpha}_t = \prod_{j=1}^t \alpha_{j/T_{\text{train}}}$ . At each training step,  $\bar{\alpha}_t$  is utilized to determine the magnitude of noise to

Hyper-parameter	Value	Hyper-parameter	Value	Hyper-parameter	Value
Scene in each Batch	8	Grasp in each Scene	64	Init LR	1e-3
LR Scheduler	Cosine	Iter	50000	Point Num	40000
Voxel side length	0.005 m	$k_{\text{trans}}$	25	$T_{\text{train}}$	1000
$T_{\text{inference}}$	200	$\beta_{\text{min}}$	0.0001	$\beta_{\text{min}}$	0.02
$\lambda_o$	1	$\lambda_g$	1	$\lambda_d$	10
$\lambda_\theta$	1	$\eta$	10		

Table 8: Hyper-parameter Setup

be added to the sample, as detailed in the main paper. At each inference step, we denoise a noisy sample  $\hat{g}_E^t$  into a less noisy sample  $\hat{g}_E^{t-dt}$  by solving the following ODE with  $t$  from 1 to 0:

$$\hat{g}_E^t - \hat{g}_E^{t-dt} = d\hat{g}_E^t = \frac{T_{\text{train}}\beta_t\sqrt{\bar{\alpha}_t}}{2\sqrt{1-\bar{\alpha}_t}}v_\Theta(\hat{g}_E^t, f_s, t)dt \quad (6)$$

Moreover, [23, 24] introduce a PDE to estimate the probability  $p(g_E|f_s)$ :

$$\frac{\partial \log p(\hat{g}_E^t|f_s)}{\partial t} = -\text{Tr}\left(\frac{\partial \bar{v}_t}{\partial \hat{g}_E^t}\right), \quad \text{where } \bar{v}_t = \frac{T_{\text{train}}\beta_t\sqrt{\bar{\alpha}_t}}{2\sqrt{1-\bar{\alpha}_t}}v_\Theta(\hat{g}_E^t, f_s, t) \quad (7)$$

Based on the above equation, we can approximate a sample’s probability  $p(g_E|f_s)$  with numerical integration during the backward process. We rank each output  $g$  of the grasp generation module using a linear combination of the estimated probability  $p(g_E|f_s)$  of the wrist pose  $g_E$  and the predicted graspness  $GS_s$  of the seed point  $s$ :

$$\text{rank}(g) = p(g_E|f_s) + \eta GS_s \quad (8)$$

**Inference Speed and Memory Cost.** Our model efficiently processes a scene point cloud comprising 40,000 points, generating 128 grasp poses and ranking them all within **0.5 seconds**. The maximum memory usage during this inference is approximately **3 GB**. These evaluations were conducted on an NVIDIA 4090 graphics card.

## 6 Implementation Details for Parallel Grippers

### 6.1 Data Filtering and Refinement

As our generative model considers all grasping poses from the dataset as successful, and since the original GraspNet-1Billion dataset [1] includes some imperfect poses, we introduce a data filtering and refinement process before training. We retain only the grasping poses with a score of  $\geq 0.9$  to ensure that all can successfully grasp the object with a friction coefficient of 0.2. To simplify motion planning, we assume that all grasps can be achieved by moving along the approaching vector and filtering out poses that would result in collisions during this movement. We also fix the depth to 4 cm and adjust the translation accordingly.

To handle poses that collide with the object and the table, we calculate the upper ( $u$ ) and lower ( $l$ ) bounds of the distance between the fingers along the original approaching vector. If the distance between any finger and the object is  $u - l < 1.5$  cm, we discard the pose. We then uniformly sample new finger positions from the adjusted lower bound  $l' = l + s$  and the adjusted upper bound  $u' = l' + \min(0.01, (u - l - 0.01) - 2s)$ , where  $s = \min(0.01, \frac{u-l-0.01}{2})$ . This ensures the fingers maintain a safe distance from the object without being too far. Finally, we calculate the intersection point of the object mesh and the new approaching vector, setting it as the seed point. Poses without a valid seed point are filtered out.

### 6.2 Graspness Definition for Gripper

For parallel grippers, after we define the intersection point as the seed point, we assign the graspness to nearby points with Eq. 3 and compute the total graspness for each point with Eq. 4, same as the dexterous hand experiments.

### 347 **6.3 Sampling Poses from Prediction**

348 Given the variability in graspness among different objects, we developed a new sampling strategy to  
349 maintain diversity and select high-quality grasping poses. First, we identify all seed points within  
350 the top 1% for graspability. For each of these seed points, we collect all points within a 2 cm radius.  
351 We then select the top 10% of these points based on graspability as new seed points and calculate  
352 grasping poses with them.

### 353 **6.4 Real-World Experiments**

354 As a lot of the objects in the LEAP Hand’s experiment are too large for our parallel gripper, we use  
355 different scenes in those two experiments as shown in Fig. 4.

## 356 **7 Additional Visualizations**

357 In Fig. 9 we present more scenes with the predictions of our network. All point clouds are colored  
358 with heatmap of model predicted graspness, with lighter color meaning higher graspness. Each  
359 scene is also dubbed with the predicted grasping pose corresponding to highest rank.

360 In Fig. 10 we show some renderings of test scenes composed of objects from ShapeNet [11].

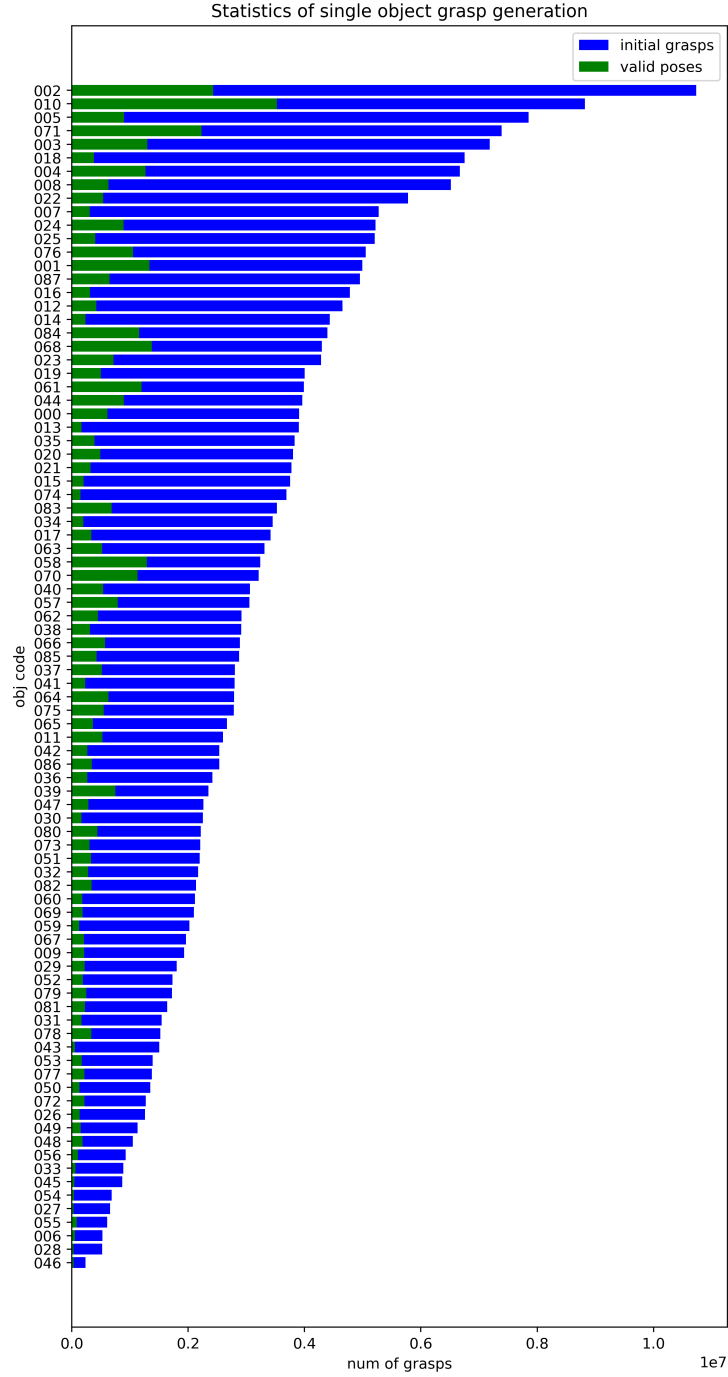


Figure 5: **Number of per-object initial grasp poses.** The proportion corresponding to valid grasps after optimization are colored green.

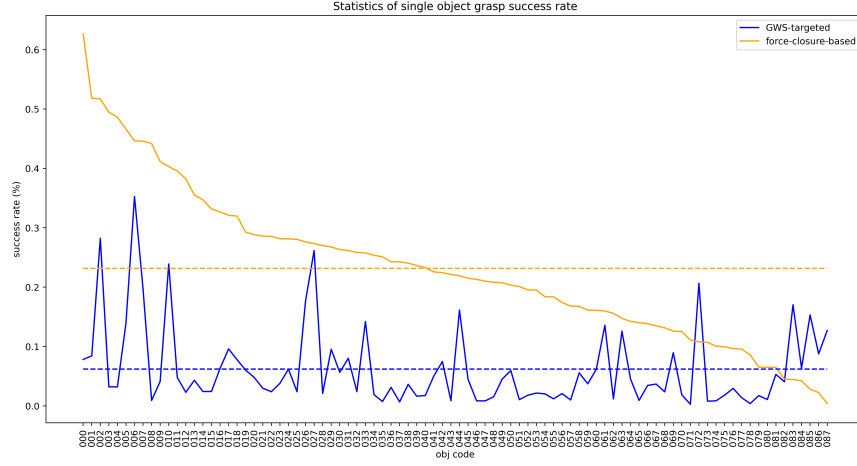


Figure 6: **Valid Rate** of single object grasp synthesis in sorted order. **Yellow** and **Blue** curves present per-object valid rates for our force-closure based optimization method (Sec.4.2.2) and GWS-based optimization method (Sec.4.2.1), respectively. Averaged success rates are drawn in dotted line, with values 24.19% and 7.91% respectively.

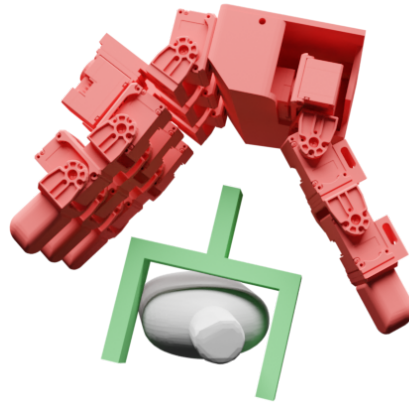


Figure 7: **Initial dexterous hand pose** superimposed with gripper grasp label at the same grasp point. We retarget gripper annotation in GraspNet-1Billion [1] to initial 6D wrist pose of dexterous hand, and use a predefined set of joint qpos for initialization.



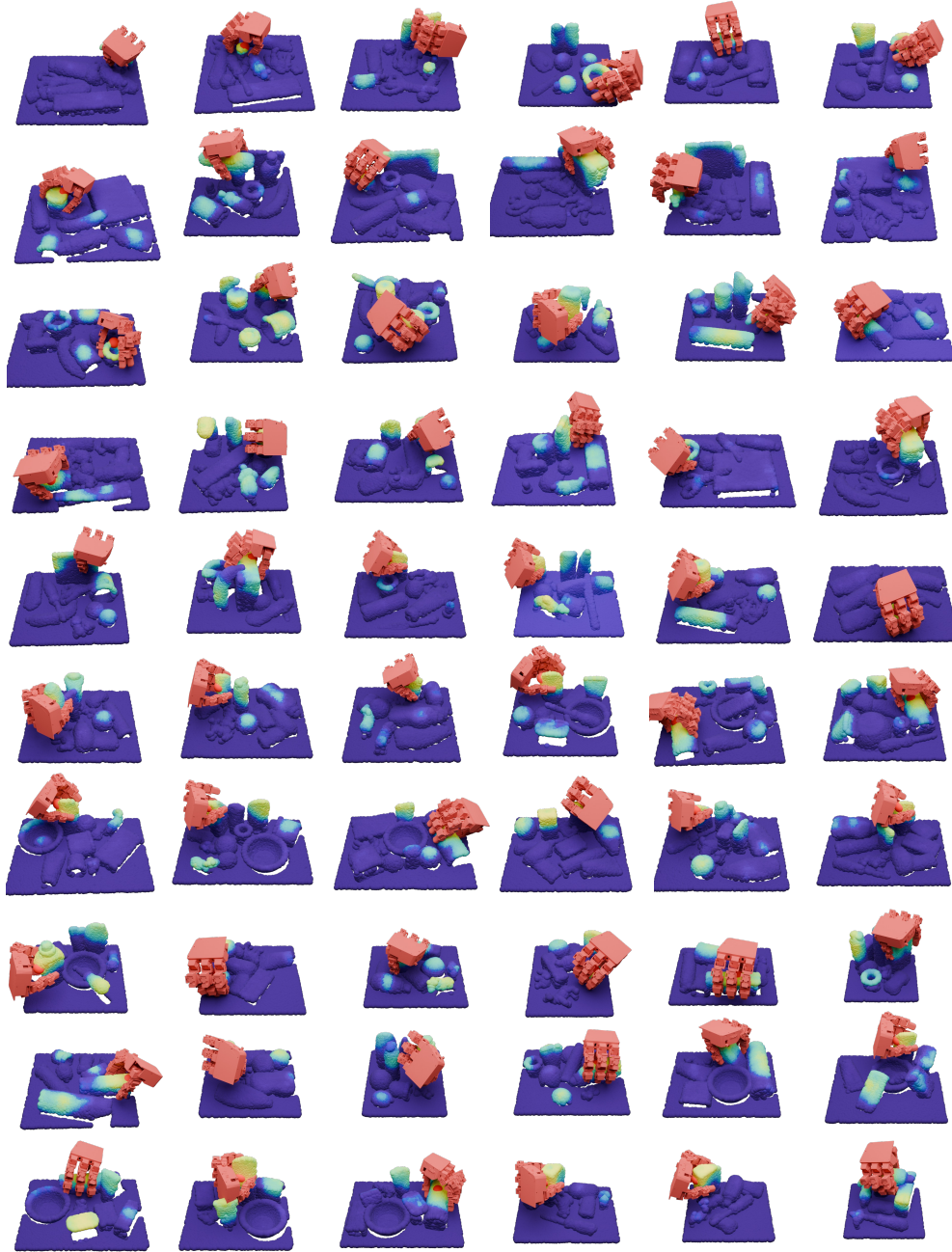


Figure 9: **Gallery visualization** of test scenes in our benchmark, corresponding to scenes 0100-0159 in GraspNet-1Billion [1]. All point clouds are colored with heatmap of model predicted graspiness, with lighter color meaning higher graspiness. Each scene is also dubbed with the predicted grasping pose corresponding to highest rank.



Figure 10: Test scenes composed of objects from ShapeNet [11].

## References

- [1] H.-S. Fang, C. Wang, M. Gou, and C. Lu. Graspnet-1billion: A large-scale benchmark for general object grasping. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11444–11453, 2020.
- [2] Y. Li, W. Wei, D. Li, P. Wang, W. Li, and J. Zhong. Hgc-net: Deep anthropomorphic hand grasping in clutter. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 714–720. IEEE, 2022.
- [3] Z. Q. Chen, K. Van Wyk, Y.-W. Chao, W. Yang, A. Mousavian, A. Gupta, and D. Fox. Learning robust real-world dexterous grasping policies via implicit shape augmentation. *arXiv preprint arXiv:2210.13638*, 2022.
- [4] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017.
- [5] H. Jiang, S. Liu, J. Wang, and X. Wang. Hand-object contact consistency reasoning for human grasps generation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11107–11116, 2021.
- [6] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017.
- [7] J. Levinson, C. Esteves, K. Chen, N. Snavely, A. Kanazawa, A. Rostamizadeh, and A. Makadia. An analysis of svd for deep rotation estimation. *arXiv preprint arXiv:2006.14616*, 2020.
- [8] J. Shi, Y. Jin, D. Li, H. Niu, Z. Jin, H. Wang, et al. Asgrasp: Generalizable transparent object reconstruction and grasping from rgb-d active stereo camera. *arXiv preprint arXiv:2405.05648*, 2024.
- [9] C. Wang, H.-S. Fang, M. Gou, H. Fang, J. Gao, and C. Lu. Graspness discovery in clutters for fast and accurate grasp detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15964–15973, 2021.
- [10] H.-S. Fang, C. Wang, H. Fang, M. Gou, J. Liu, H. Yan, W. Liu, Y. Xie, and C. Lu. Anygrasp: Robust and efficient grasp perception in spatial and temporal domains. *IEEE Transactions on Robotics*, 2023.
- [11] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [12] Q. Dai, J. Zhang, Q. Li, T. Wu, H. Dong, Z. Liu, P. Tan, and H. Wang. Domain randomization-enhanced depth simulation and restoration for perceiving and grasping specular and transparent objects. In *European Conference on Computer Vision*, pages 374–391. Springer, 2022.
- [13] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. <http://pybullet.org>, 2016–2021.
- [14] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
- [15] J. Chen, Y. Chen, J. Zhang, and H. Wang. Task-oriented dexterous grasp synthesis via differentiable grasp wrench boundary estimator. *arXiv preprint arXiv:2309.13586*, 2023.
- [16] 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.

- 406 [17] B. Sundaralingam, S. K. S. Hari, A. Fishman, C. Garrett, K. Van Wyk, V. Blukis, A. Millane,  
407 H. Oleynikova, A. Handa, F. Ramos, et al. curobo: Parallelized collision-free minimum-jerk  
408 robot motion generation. *arXiv preprint arXiv:2310.17274*, 2023.
- 409 [18] H. Dai, A. Majumdar, and R. Tedrake. Synthesis and optimization of force closure grasps via  
410 sequential semidefinite programming. *Robotics Research: Volume 1*, pages 285–305, 2018.
- 411 [19] C. Choy, J. Gwak, and S. Savarese. 4d spatio-temporal convnets: Minkowski convolutional  
412 neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern  
413 recognition*, pages 3075–3084, 2019.
- 414 [20] D. Misra. Mish: A self regularized non-monotonic activation function. *arXiv preprint  
415 arXiv:1908.08681*, 2019.
- 416 [21] C. Durkan, A. Bekasov, I. Murray, and G. Papamakarios. nflows: normalizing flows in Py-  
417 Torch, Nov. 2020. URL <https://doi.org/10.5281/zenodo.4296287>.
- 418 [22] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural  
419 information processing systems*, 33:6840–6851, 2020.
- 420 [23] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural ordinary differential  
421 equations. *Advances in neural information processing systems*, 31, 2018.
- 422 [24] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based gen-  
423 erative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*,  
424 2020.