

ANYDEXGRASP: LEARNING GENERAL DEXTEROUS GRASPING FOR ANY HANDS WITH HUMAN-LEVEL LEARNING EFFICIENCY

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce an efficient approach for learning dexterous grasping with minimal data, advancing robotic manipulation capabilities across different robotic hands. Unlike traditional methods that require millions of grasp labels for each robotic hand, our method achieves high performance with human-level learning efficiency: only hundreds of grasp attempts on 40 training objects. The approach separates the grasping process into two stages: first, a universal model maps scene geometry to intermediate contact-centric grasp representations, independent of specific robotic hands. Next, a unique grasp decision model is trained for each robotic hand through real-world trial and error, translating these representations into final grasp poses. Our results show a grasp success rate of 75-95% across three different robotic hands in real-world cluttered environments with over 150 novel objects, improving to 80-98% with increased training objects. This adaptable method demonstrates promising applications for humanoid robots, prosthetics, and other domains requiring robust, versatile robotic manipulation.

1 INTRODUCTION

Grasping, as a fundamental problem of prehensile manipulation, holds significant importance in robotics. Over the past decades, diverse mechanical structures for robotic hands have been developed. Visually guided dexterous grasping is in high demand to enable robots to interact effectively with their environments. This ability also plays a crucial role in the context of intelligence. Throughout human evolution, early humans developed the capability for precise grip Skinner et al. (2015), which enabled tool use and is believed to have facilitated the evolution of the human species Almeija et al. (2010); Kivell (2015). From the perspectives of both advancing robotics and promoting embodied intelligence, it is essential to design a learning framework that efficiently equips different robotic hands with visually guided dexterous grasping capabilities.

To make such a grasping system practically useful, it should use a single commodity camera, observe environments with cluttered objects, handle perception noise and generate a set of dexterous grasp poses that can be selected by subsequent tasks. Due to the challenges of the problem, early research focused on generating dexterous grasp poses given a single, complete object mesh, utilizing either analytical Miller & Allen (2004); Rosales et al. (2011); Liu et al. (2021; 2020) or learning-based approaches Li et al. (2023). The idea is to decouple the grasping system into 6D pose estimation and grasp poses generation based on the object CAD model. However, the requirement for the object mesh limits its ability to handle new object shapes.

It is challenging to detect grasp poses for unseen objects based on partial-view perception. Some recent methodologies pursue mesh completion using partial point clouds Lundell et al. (2021); Wei et al. (2022; 2024), followed by grasps generation on the complete mesh. However, the error introduced by perception noise and mesh completion often results in inaccurate grasp analysis. An increasing amount of research has attempted to learn the mapping from raw partial observation to grasp poses within a single network. Due to the highly nonlinear property of this mapping, extensive training data is required. Two data sources are commonly adopted: human grasping demonstrations Gupta et al. (2016); Christen et al. (2019); Qin et al. (2022); Mandikal & Grauman (2022); Wei et al. (2024); Shaw et al. (2024) or data from simulated environments Brahmhatt et al. (2019);

054 Corona et al. (2020); Grady et al. (2021); Li et al. (2023); Wang et al. (2023). However, both meth-
055 ods have their limitations. The former approach struggles to accurately capture hand gestures and is
056 confined to robotic hands resembling human anatomy. The latter requires substantial effort to build
057 the simulation environment, annotate grasp poses, and transfer algorithms from simulation to the
058 real world. These challenges limit current grasping systems to simple scenarios, typically involving
059 a single object at a time from a limited set. No prior work demonstrates robust grasping in cluttered
060 environments from partial-view perception in the real world.

061 Most critically, even if these challenges are overcome, the policy obtained with substantial efforts
062 is only suitable for a specific robotic hand each time. The end-to-end learning paradigm implic-
063 itly encodes the information about hand kinematic structure, relevant state information and grasp
064 quality in the weights, making it difficult for models to share computation between different hands.
065 Consequently, we need to repeat the tedious data generation and policy training pipeline for each
066 hand.

067 We identified two main bottlenecks for efficiently learning visually guided dexterous grasping for
068 different robotic hands: the requirements for extensive training data for each hand, and the inability
069 to share computation across different hands. These bottlenecks arise from attempting to learn the
070 mapping from raw observation to grasp poses with an end-to-end network. In this paper, we revisit
071 this paradigm. We hypothesize that if there exists a low-dimensional intermediate state space that
072 encapsulates grasp information, then the mapping from this state space to grasp poses can be learned
073 more efficiently than the original mapping, requiring less training data. Moreover, if such a state
074 space is transferable across different robotic hands, it could be shared without the need to retrain a
075 state estimator each time. Note that such a state space should not require object knowledge during
076 inference, in order to generalize to unseen objects.

077 Recognizing this potential, we aimed to identify such a state space. For the grasping problem, the
078 robot needs to decide its grasp forces on each finger, based on the grasp matrix, surface normals of
079 contact points and friction coefficient Dai et al. (2018). From visual perception, the information we
080 can extract is the positions and normals of potential contact points, where positions are linked to the
081 grasp matrix and normals determine the orientations of friction cones. Based on this observation, we
082 introduce a novel intermediate representation for multi-finger grasping, referred to as the Contact-
083 centric Grasp Representation (CGR), which encapsulates contact information on the object’s surface
084 and possesses SE(3)-equivalent property.

085 Based on this representation, we present AnyDexGrasp, a novel methodology that can effectively
086 learn dexterous grasping for different hands on a modest set of training objects. In this method,
087 the multi-finger grasp detection problem that maps raw perception to grasp poses is divided into
088 two steps. In the first step, we train a general representation model that maps single-view partial
089 observations to contact-centric grasp representations. A large-scale dataset is annotated to train
090 this model. After training, it can be applied to different hands without fine-tuning. In the second
091 step, we map the contact-centric grasp representations to a set of grasp proposals through a hand-
092 specific mapping, and then learn a hand-specific classifier to evaluate each grasp proposal. This
093 classifier takes a contact-centric grasp representation and a grasp proposal as input and maps them
094 to the probability of grasp success. The training data is collected by real-world trial and error. We
095 empirically observed that this mapping is significantly easier to learn, requiring merely hundreds
096 of trial-and-error attempts. It dramatically reduces the cost of real-world learning and allows our
approach to work for different types of robotic hands efficiently.

097 We evaluate the effectiveness of our method using three different robotic hands, each featuring
098 three to five fingers. Our system is first trained on 144 objects, with approximately 2,000 to 8,000
099 grasp attempts, depending on the robotic hand. On a diverse set of 150 previously unseen objects,
100 including deformable and adversarial items, our approach achieves an average grasp success rate
101 ranging from 80% to 98% across different hands. Notably, this performance is achieved in cluttered
102 scenarios, demonstrating the effectiveness of our approach.

103 In addition we explore further reductions in training samples required for our grasp learning
104 paradigm. We limit the training objects to 40 and reduce the grasp attempts to approximately 400
105 to 1,000 depending on the robotic hand. Even with this limited amount of training data, our system
106 consistently achieves grasp success rates ranging from 75% to 95% during real-world testing. No-
107 tably, our experiments also highlight the potential for further reductions of training samples, with the

108 ability to decrease the training object number to 30 and the total grasp attempts to 200 for a three-
109 finger hand, without decreasing the grasp performance by a large margin. Such learning efficiency
110 allows robots to master visually guided grasping in a matter of hours in the real world, surpassing
111 the learning efficiency of human infants.

112 We conduct a series of analyses to clarify why our two step learning method is so efficient. In the first
113 step, we perform a geometry coverage analysis, showcasing that by scaling up data in the correct
114 dimension, the local geometries on just 40 objects can effectively cover a wide range of unseen
115 objects. This explains the generalization capabilities with a small number of training objects. In
116 the second step, we provide various perspectives illustrating how our proposed contact-centric grasp
117 representation serves as a robust state space for grasp decision, which allows the model to learn from
118 just hundreds of real-world trial-and-error attempts.

119 This paper represents a significant step toward the efficient realization of dexterous robotic grasping,
120 with the potential to revolutionize various applications, from advanced humanoid robots to prosthetic
121 hands.

123 2 RESULTS

124
125 We first learn a general hand-agnostic representation model based on an offline annotated, large-
126 scale dataset. Once the representation model is learned, we use the predicted contact-centric grasp
127 representation as a new state space for the problem of multi-finger grasping. For three different
128 robotic hands (DH-3, Allegro and Inspire), we learn grasping in the real world directly through trial
129 and error. We discretize each robotic hand’s configuration space into several predefined grasp types,
130 following human grasping taxonomy Cutkosky et al. (1989). We start with thousands of trial-and-
131 error grasp attempts, and gradually reduce the number to hundreds of attempts in later experiments
132 to demonstrate the efficiency of our learning paradigm. We also vary the number of training objects
133 from 144 to 40 and even 30, to verify the generalization ability of our grasp system.

134 2.1 LEARNING DEXTEROUS GRASPING WITH 144 OBJECTS

135
136 We first employ a training set of 144 objects to train the grasp decision model. Approximately 1,000
137 grasp samples are collected for each grasp type of the DH-3 and Inspire Hand, and 200 grasp samples
138 are collected for each grasp type of the Allegro Hand, forming the basis for our learning process.
139 The amount of training objects and grasp samples would be gradually reduced in later sections to
140 verify the effectiveness of our method.

141 2.1.1 DEXTEROUS GRASPING ON DAILY OBJECTS

142
143 We systematically evaluate the success rates of our approach on testing objects from the first five
144 categories commonly encountered in our daily activities. The average success rates achieved by the
145 three distinct robotic hands are 97%, 78%, and 83%, respectively. Movies S1, S2 and S3 record
146 the grasping process. In contrast, the success rates of the baseline method using heuristic sampling
147 and collision detection reach only 66%, 51%, and 58%. A detailed breakdown of success rates for
148 each object category is presented in Figure 2A. Compared to the baseline method, the substantial
149 improvements across this extensive test set demonstrate the effectiveness of our proposed represen-
150 tation and approach.

151
152 Several noteworthy points are hereby highlighted. Firstly, the 3-finger gripper attains an average
153 success rate of 97% across over 100 real-world objects, surpassing even the performance of previous
154 state-of-the-art parallel-gripper algorithm Fang et al. (2023b). Secondly, for deformable objects
155 within the textile and food categories, the grasp success rates across different grippers show no
156 significant degradation. In some cases, they even slightly outperform the overall success rate, despite
157 the absence of explicit training on deformable objects. This observation emphasizes the remarkable
158 generalization capacity of data-driven methods. We observe that deformable objects tend to comply
159 with the gripper during the grasping process, making them easier to be successfully grasped.

160 Regarding grasping speed, our system takes an average of 0.5 second to generate 200 grasp poses in
161 a cluttered scene. Additionally, an extra collision detection step utilizing scene partial point cloud
and hand mesh is performed. It takes 20 seconds on our CPU using Open3D library Zhou et al.

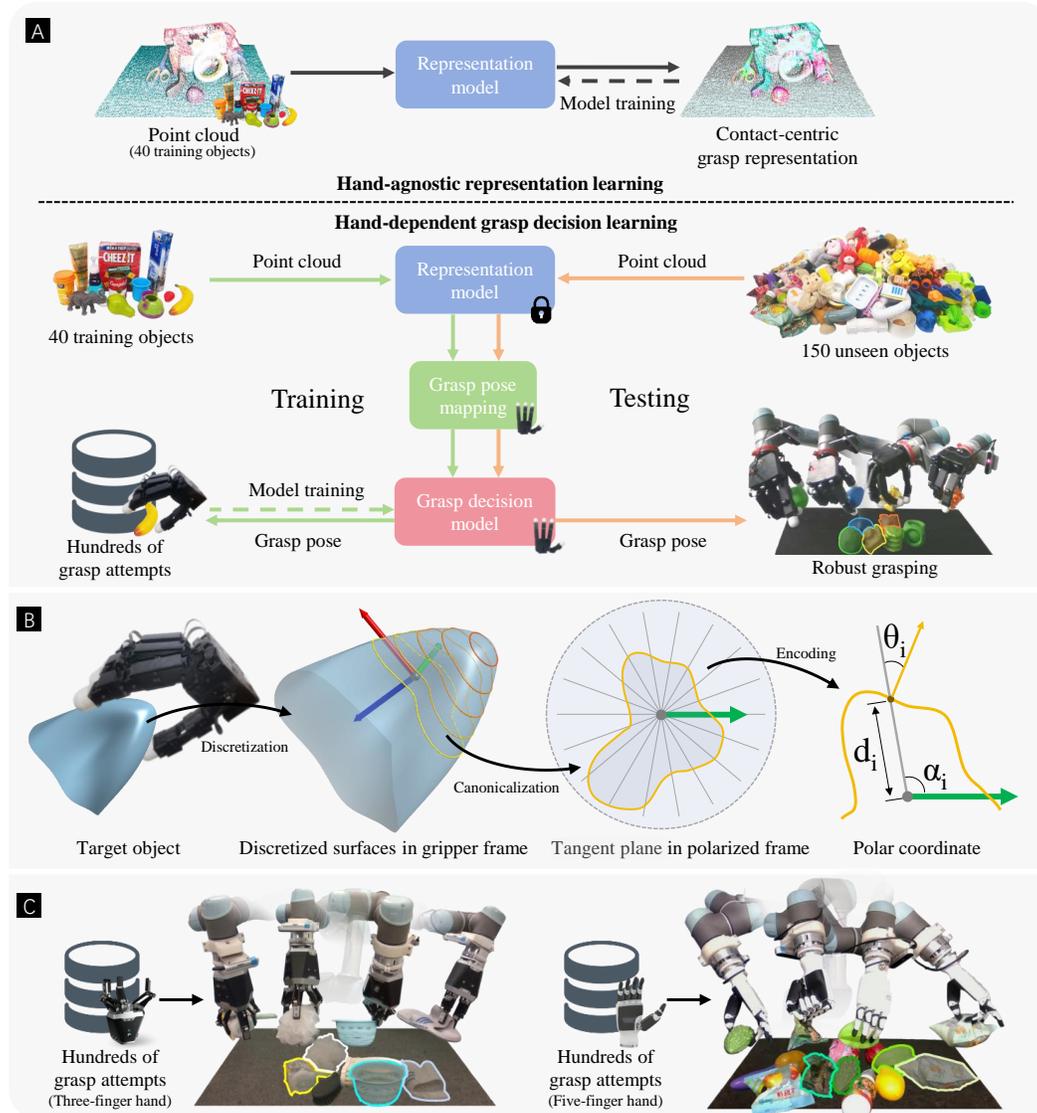


Figure 1: **The overview of our method.** (A): Our method consists of two steps. The first step is to train a representation model on partial-view point cloud. The training set only consists of 40 objects. The second step would fix the representation model, and train a grasp decision model that takes the grasp-centric contact representation as input and outputs the grasp success score, based on hundreds of real-world trial-and-error attempts. The grasp algorithm is tested thoroughly on hundreds of unseen objects. (B): Illustration of contact-centric grasp representation. A local geometry is discrete into several tangent planes along the approach direction of a robotic hand. Each tangent surface is transformed into the polarized coordinate frame of the robotic hand. The shape of the surface is encoded into discretized points and normal representation in the polar coordinate. (C): Our experiments are also carried out on a three-finger hand and a five-finger hand and demonstrate excellent performance.

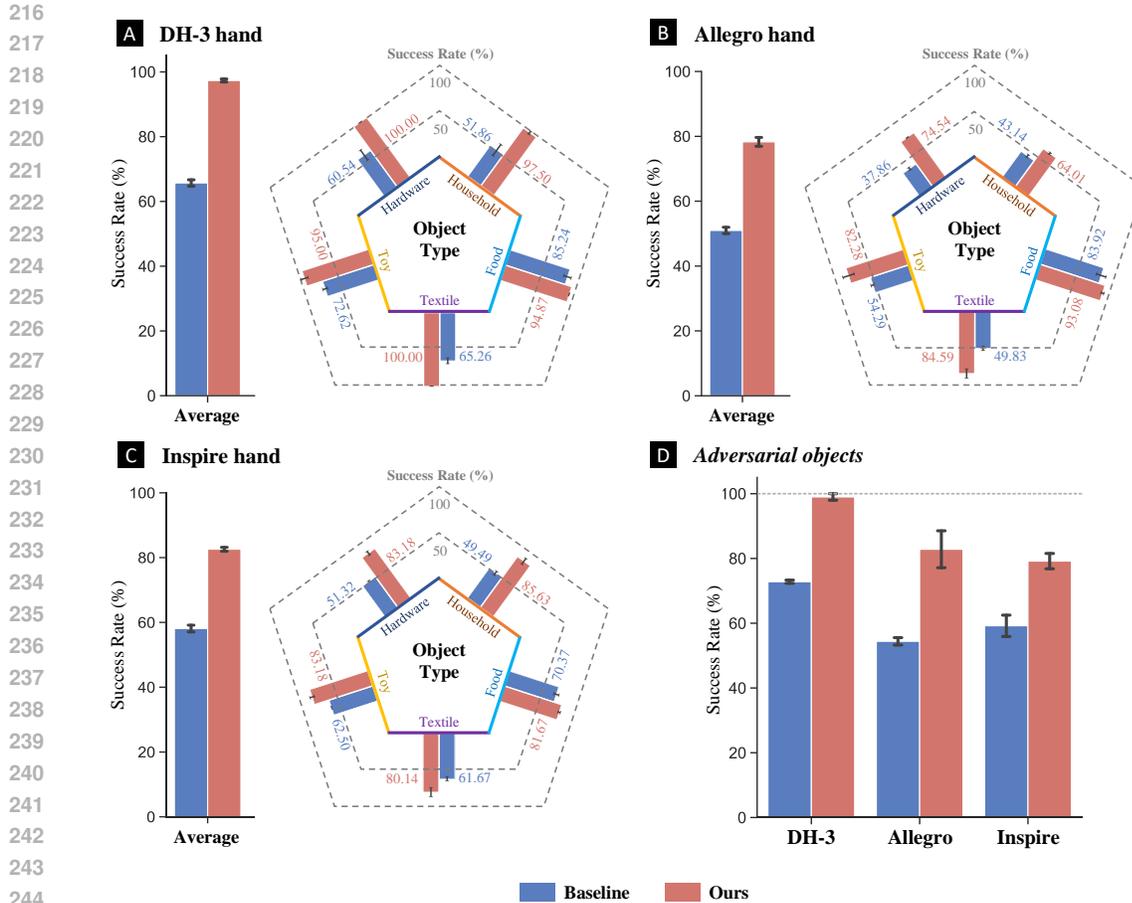


Figure 2: **Success rates on the testing set after training on abundant real-world data.** (A): The averaged and detailed success rates of the DH-3 hand on five object categories commonly encountered in our daily activities. (B): The averaged and detailed success rates of the Allegro hand. (C): The averaged and detailed success rates of the Inspire hand. (D): The success rates on the adversarial objects of three robotic hands.

(2018). Although this step could be accelerated through advanced collision detection technology or hardware acceleration, this aspect falls outside the scope of this paper.

2.1.2 DEXTEROUS GRASPING ON ADVERSARIAL OBJECTS

In addition to daily objects, we extend our method’s evaluation to more challenging adversarial objects. These objects encompass 13 human-selected items from DexNet Mahler et al. (2017) and 49 program-generated objects from EGAD! evaluation set Morrison et al. (2020), characterized by distinct shapes and varying grasp difficulties. Prior literature shows a performance degradation of parallel grasping on adversarial versus daily objects Fang et al. (2023b). To the best of our knowledge, this is the first comprehensive evaluation of a multi-finger grasping algorithm on adversarial objects in real-world scenarios.

Success rates for the three distinct robotic hands are reported in Figure 2B, where our system achieves 99%, 82%, and 79% success rates, respectively. Movies S4, S5 and S6 record the grasping process. In contrast, the baseline method achieves success rates of 72%, 54%, and 59%. Remarkably, the performance on adversarial objects is on par with daily objects for all the robotic hands, highlighting the promising generalization ability of our dexterous grasping system. It surprises us since previous results from parallel grippers show a dramatic performance degradation. We presume

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

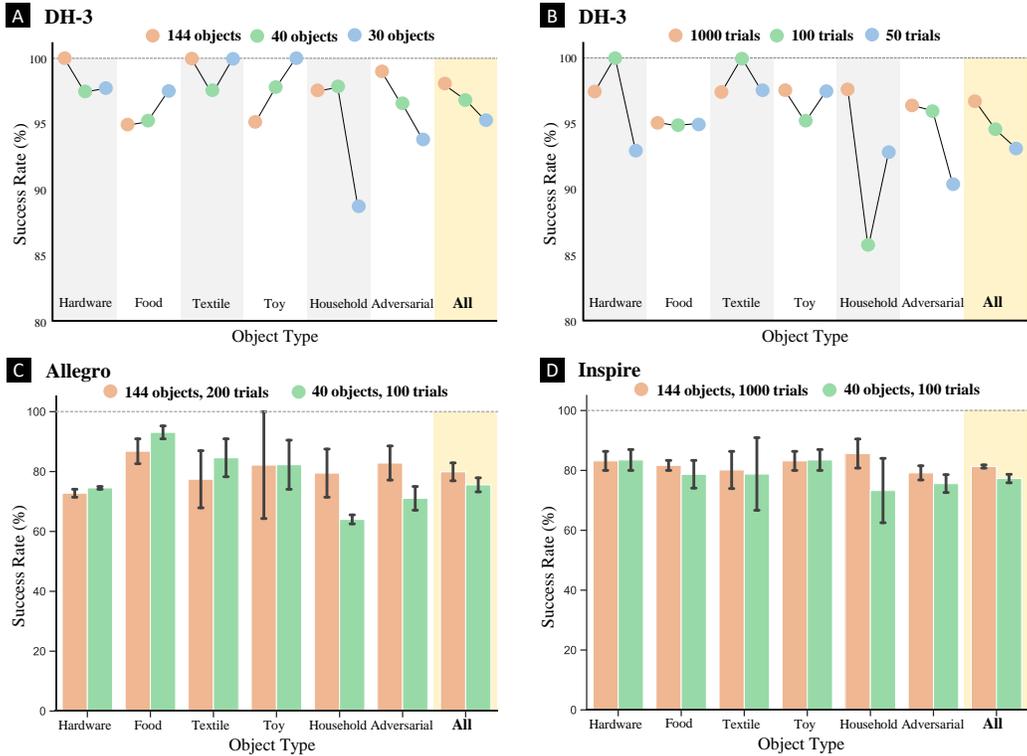


Figure 3: **Success rates on the testing set after training on reduced real-world data.** (A): We reduce the training object number from 144 to 40 and 30 respectively and test the success rates on different categories of objects. (B): With 40 training objects, we reduce the data from around 1000 trials per grasp type to 100 trials and 50 trials respectively. (C) and (D): When reducing the training data on fewer training objects and fewer grasp attempts, success rates on both Allegro hand and Inspire hand only decrease slightly, showing good generalization ability and high learning efficiency of our method.

that the additional fingers can improve the grasping ability, and the adversarial objects designed for parallel grippers do not pose significant challenges in multi-finger cases.

In the subsequent sections, unless otherwise stated, we proceed to conduct experiments on all 150 objects including daily ones and adversarial ones.

2.2 REDUCING REAL-WORLD TRAINING BURDEN

In previous experiments, we collected a considerable volume of real-world training data, approximately 1,000 trials per grasp type, on the full set of 144 training objects. In this section, we aim to alleviate the demands of real-world training by assessing the model’s performance under reduced object and trial conditions. For the sake of simplicity, our evaluation concentrates on the 3-finger hand in this section.

We initiate this exploration by reducing the number of objects utilized in real-world trial-and-error attempts. Two smaller sets consisting of 40 and 30 objects were adopted for our experiments. The details of the training object sets are given in Materials and Method. To ensure a fair comparison with the original object set, we collect an equivalent number of around 1,000 grasp samples for each grasp type, ensuring the convergence of the grasp decision model.

Figure 3A presents the experimental results. As the training object count reduces from 144 to 40, the grasp success rate achieves 96.7% on all testing objects, representing only a marginal decrease of 1.1%. Such subtle performance degradation, given a nearly 3/4 reduction in training objects, showcases the robustness of our approach. A further reduction to 30 training objects results in an

324 overall grasp success of 95.1%. This translates to a further decrease of 1.6%, yet the performance
325 remains notably promising. Our experiments demonstrate that, with proper learning methods, the
326 thousands of training objects adopted in previous systems Mahler et al. (2019); Wang et al. (2023)
327 are not necessary.

328 Since the performance degradation when reducing the object set from 144 to 40 and from 40 to 30
329 is comparable, while reducing the training object set from 40 to 30 does not significantly lower the
330 training burden, we opt to proceed with the 40-object training set for subsequent experiments.

331 We then explore the impact of reducing the number of trials and errors for each grasp type. On the
332 40 training objects, we reduce trials and errors from approximately 1,000 attempts per grasp type to
333 100 attempts and 50 attempts, respectively. Figure 3B presents the real robot testing results. When
334 training with 100 trials per grasp type, the success rate reaches 94.5% on average on all objects. We
335 show the whole grasping process in Movie S7. This success rate is strikingly high given the limited
336 number of real-world training samples. Previous literature Xu et al. (2023); Liu et al. (2023) often
337 required millions of grasp attempts in simulation to achieve grasping proficiency. Further reducing
338 the trials to 50 attempts per grasp type yields a success rate of 93.1% on all objects. These results
339 demonstrate the high learning efficiency of our method, which requires only a small number of grasp
340 attempts for convergence. In our following experiments, given the already high efficiency of 100
341 trials per grasp type, we adopt this setting for learning.

342 2.3 DEXTEROUS GRASP LEARNING WITH 40 OBJECTS AND 100 ATTEMPTS

343 In the previous section, we demonstrated the robust grasping policy acquired by the 3-finger gripper
344 through a significantly limited amount of training data and real-world attempts. In this section, we
345 extend the validation of such a learning paradigm to the other two robotic hands utilized in this
346 study.
347

348 We directly assess the performance of training using 40 objects with 100 trials for each grasp type
349 on the four-finger and five-finger hands. Depending on the number of grasp types, the total real-
350 world training samples amount to 1,000 and 800 for these two robotic hands, respectively. This
351 significantly reduced volume of real-world training samples, nearly 1/10 of the original experiments,
352 presents a territory in grasp learning that is unexplored by previous work.

353 Figure 3C and Figure 3D display the detailed success rates of real-world experiments. The average
354 success rates stand at 75%, and 77% for all objects. The grasping process is recorded in Movies
355 S8, S9, S10 and S11. It's striking that the success rates show minimal decreases compared to the
356 original performance. This observation demonstrates the substantial learning efficiency enabled
357 by our methodology. Such proficiency allows diverse robotic hands to acquire dexterous grasping
358 ability in real-world settings.

359 Notably, this efficiency surpasses that observed in human infants, who typically require months of
360 practice to develop visually guided grasping skills. The grasp success rates for human infants reach
361 61.9% at 8 months old Domellöf et al. (2015), which involves thousands of practice attempts starting
362 at 4 months old Newell et al. (1989). It is noteworthy that our grasping results are achieved based
363 solely on visual perception, with no tactile feedback.

364 2.4 INFLUENCE OF GRASP TYPES

365 2.4.1 ACCURACY OF DIFFERENT GRASP POSES

366 In the above experiments, we have shown that our method can enable efficient grasp learning with
367 high success rates. Here we further analyze the success rates of each robotic hand with a detailed
368 breakdown according to their respective grasp types. For clarity purposes, we number each grasp
369 type, as illustrated in Figure 8. The results trained on 40 objects and 100 grasp attempts per grasp
370 type are adopted for analysis, as depicted in Figure 4A. For each hand, we can see that different grasp
371 types have different difficulties in dealing with grasping. Usually, the success rates after learning
372 are dramatically higher than the baseline method. However, there also exist some exceptions. For
373 example, the grasp type 1 of the five-finger Inspire hand after training yields a close success rate
374 to the baseline. After inspection, we found that this grasp type was selected fewer times after the
375
376
377

378 training. We anticipated that other grasp types might be more confident to grasp if the objects can
379 be grasped by multiple types, which leaves some hard cases for this grasp type.
380

381 382 2.4.2 DISTRIBUTION OF GRASP TYPES 383

384 A natural question that arises is whether the system learned by our method can demonstrate a variety
385 of grasp types. From the example above, it is possible that the system may achieve high success
386 rates by favoring one or two grasp types while ignoring diversity. To address this, we analyze how
387 frequently each grasp type is selected during testing to verify whether our system indeed learns
388 diverse grasp poses. To quantify the frequency of each grasp type, we normalize by dividing the
389 number of grasp attempts for each type by the total number of grasp attempts across all types.
390

391 To establish a baseline, we examine the frequency of grasp types obtained by the baseline method,
392 reflecting the inherent frequency determined solely by collision detection. Grasp types prone to
393 collision with the scene naturally constitute a smaller fraction among all types. This baseline grasp
394 type frequency serves as a reference for natural distribution. The top row of Figure 4B illustrates the
395 grasp type frequency for different robotic hands. Notably, the three-finger hand exhibits a balanced
396 distribution, whereas the four- and five-finger hands display more unbalanced distributions. This
397 discrepancy arises from the fact that the fingertips of the three-finger hand consistently point in the
398 same direction along the approach vector, resulting in similar collision situations across different
399 types. Conversely, the four- and five-finger hands exhibit types with greater variance, including
some that are prone to collision with the scene.

400 Then, we present the frequency of grasp pose after employing our learned system. The statistics are
401 given in the second row of Figure 4B. For the three-finger hand, type 3 presents an increasing ratio
402 among all grasp types. The reason is that this grasp type presents a higher success rate, and usually
403 has a higher grasp quality score than other grasp types. However, the other three grasp types are also
404 frequently selected. For the four- and five-finger hand, the grasp frequency is similar to the baseline
405 method. These results affirm that our learned system adeptly captures diverse grasp poses, achieving
406 high success rates without compromising grasp diversity.
407

408 409 2.4.3 REDUCING GRASP TYPES 410

411 Another question for multi-finger grasping is whether employing multiple grasp types is necessary,
412 given the argument that a single power grasp might be sufficient for good results. However, we
413 argue that incorporating multiple types enhances flexibility, particularly when faced with cluttered
414 scenarios. To prove that, we conducted a targeted experiment to compare grasping outcomes with
415 varying numbers of grasp types. Specifically, we employed the best grasp model trained with 144
416 objects for the Allegro hand, initially defined with 10 grasp types. In our experiment, we compared
417 the original model with two modified versions that use fewer grasp types. The first version was
418 limited to a single grasp type, specifically the one that achieved the highest overall success rate
419 across all types. The second version used a subset of the five most effective grasp types, chosen based
420 on their individual success rates. For simplicity, the evaluation focused exclusively on adversarial
421 objects due to the performance similarity with that on the entire object set. The resulting success
422 rates are detailed in Figure 4C.

423 The original method, employing 10 grasp types, achieved an 80.3% success rate on the test set. In
424 contrast, utilizing only a single grasp type led to a reduction in the grasp success rate to 67.3%.
425 Employing five grasp types performed better, resulting in a success rate of 77.6%, but is still inferior
426 to the original method. Our experimental results show that increasing the number of grasp types
427 can improve overall grasp success rates. One reason for this improvement is that a greater variety
428 of grasp types provides more flexibility, enabling the hand to better adapt to different object shapes,
429 sizes, and orientations. Additionally, using multiple grasp types can increase tolerance for collisions,
430 allowing the hand to adjust its grasping strategy based on spatial constraints, particularly in cluttered
431 environments. It is noteworthy that, on the other hand, our results also reveal that the benefits derived
from further adding grasp types would eventually saturate. Thus, it is reasonable to adopt diverse
yet limited grasp types, which optimizes both grasp success rates and learning efficiency.

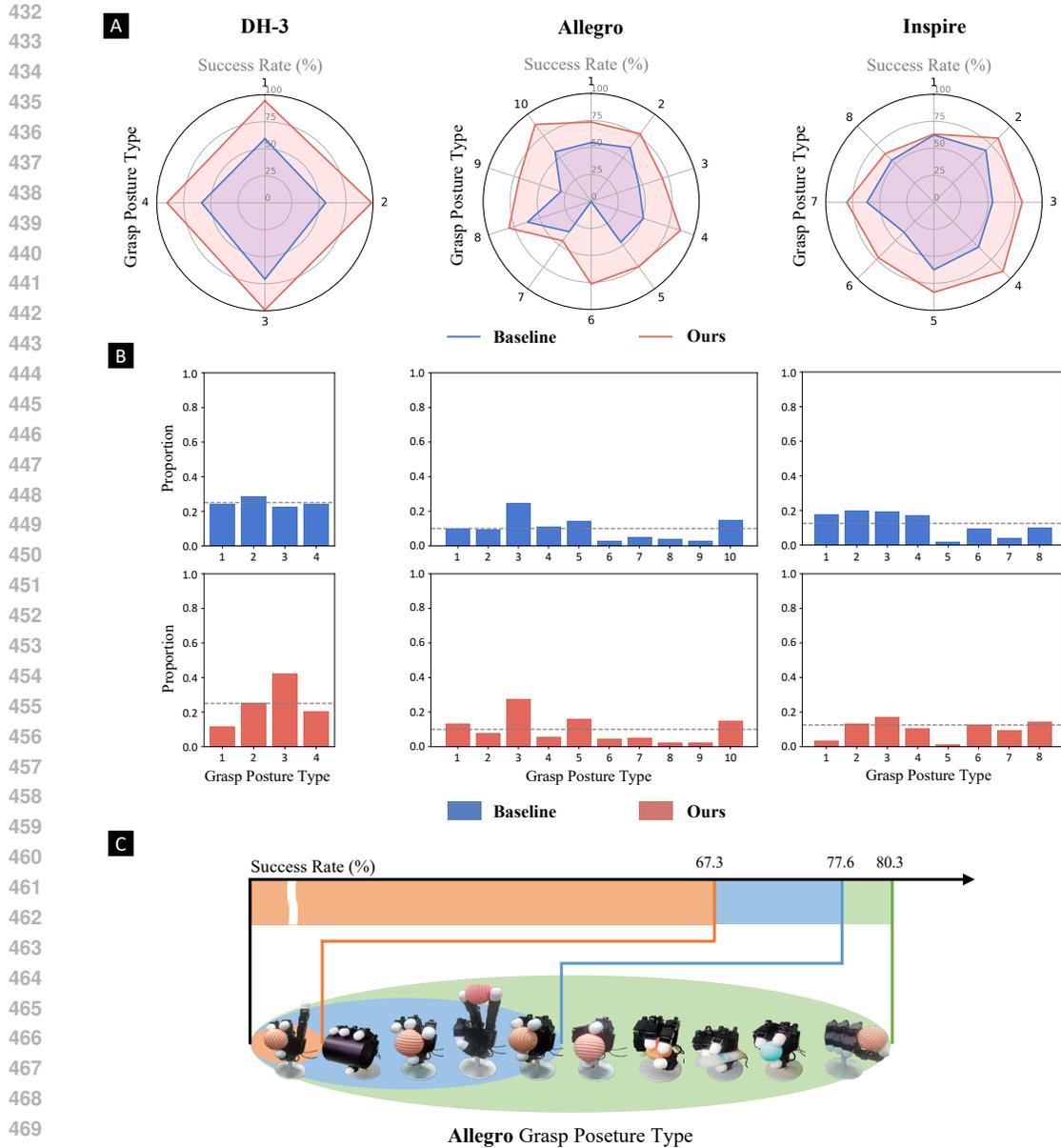


Figure 4: **Analysis of the influence by grasp type.** (A) A breakdown analysis of grasp success rates on different grasp types for each robotic hand. (B) The selected frequency of different grasp types for each robotic hand during testing. (C) Grasping success rates when using different portions of grasp types for the allegro hand.

3 DISCUSSION

The evolution of visually guided dexterous grasping methodologies within robotics has developed two prominent paradigms: the 6D pose estimation paradigm and the end-to-end grasp learning paradigm. The former relies on the precise estimation of an object’s 6D pose and then calculates the hand pose accordingly. It can transfer across different robotic hands easily, but requires prior knowledge of the object’s model. On the other hand, the end-to-end grasp learning models do not require explicit object knowledge, yet the trained models lack transferability across different robotic hands.

Our proposed approach explores a middle ground between these two paradigms, which combining the advantages of both. By developing a contact-centric grasp representation that encapsulates the scene’s contact information, we eliminate the need for an object’s model beforehand. The CGR preserves critical information pertinent to grasp quality, endowing our system with adaptability and applicability across different morphologies of robotic hands. Moreover, by eliminating the need for an accurate kinematic model, which was frequently used in previous work learned in simulation Xu et al. (2023); Wan et al. (2023), our method is suitable for soft hand grasp learning.

It is surprising to see that our grasp system can be learned for different hands so efficiently. Previous work for multi-finger grasping usually require thousands of objects and millions of grasp samples Mahler et al. (2019); Eppner et al. (2021). And in the deep learning era, it seems to be an underlying rule that we need to train a robot system on as many objects as possible to have good generalization ability. However, the satisfactory performance of our system breaks this intuition. What are the key factors for our method to learn such efficiently and generalize so well? There are two aspects of learning efficiency in our system, the first is that we only need 40 objects and the second is that we only need hundreds of trials for each hand. Here we discuss how our method achieves efficiency in these two aspects.

Representation plays a crucial role in our system for real-world learning, as it must map different geometries into a contact-centric grasp representation. How can the system, trained on only 40 objects, generalize to hundreds of unseen objects? We address this by conducting a geometry coverage analysis, revealing that **scaling up data along the right dimension** is key to improving the model’s generalization ability.

Our representation model takes a scene point cloud as input and outputs contact-centric grasp representations (CGRs). To train this model, we need a dataset containing scene point clouds with annotated CGRs across various geometries. Since CGRs depend on local geometry, a representative dataset must include diverse local geometries to effectively train deep networks. While many researchers intuitively attempt to collect more training objects to achieve this, our geometry coverage analysis demonstrates that more objects do not necessarily lead to richer local geometries.

We begin by defining the local geometry used in our analysis. Specifically, the representation network operates in a partial observation scenario, where it infers contact positions and normals on unobserved surfaces based on the observed geometry. Although each normal and contact point is predicted independently, we consider the minimal continuous components of local patches, enclosed by the simplest form of grasping—an antipodal grasp—as a foundational element in our analysis for consistency.

Next, we define geometry coverage in the analysis. Given a training and testing object set, a local geometry on a test object is considered “covered” if it closely resembles a local geometry from the training object set. We define similarity by a chamfer distance smaller than 1mm, with examples illustrated in Figure 5A. For a training dataset, we can assess the diversity of local geometries by counting the number of covered geometries on test objects.

In practice, when constructing a training dataset, we need to generate labels for a fixed number of local geometries selected from the training objects, constrained by computational resources. There are two possible dimensions along which to collect more local geometries: increase the number of training objects or increase the sample density on each training object. To assess which dimension is more effective for increasing local geometry diversity, our analysis is conducted as follows. We collect two training object sets: \mathbb{S} , with 40 objects, and \mathbb{L} , with 144 objects. For each object set, we sample 1 million and 4 million local geometries from each object on average, respectively (sampling details given in supplementary material). This combination results in four different training datasets. The testing object set for the coverage analysis is the EGAD! test set, which contains over 2000 complex, program-generated objects. We sample around 400 local geometries on each test object and evaluate if they are covered. Figure 5B shows the number of covered geometries on each test object, considering different training sets and sampling densities.

Surprisingly, we found that increasing the number of training objects does not significantly increase the coverage rate on the test set. However, increasing the sampling density of local geometries per object leads to a dramatic increase in coverage—even when the total number of sampled geometries is similar to increasing the object count. This result demonstrates that increasing sample density for each training object is far more impactful than increasing the number of objects.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

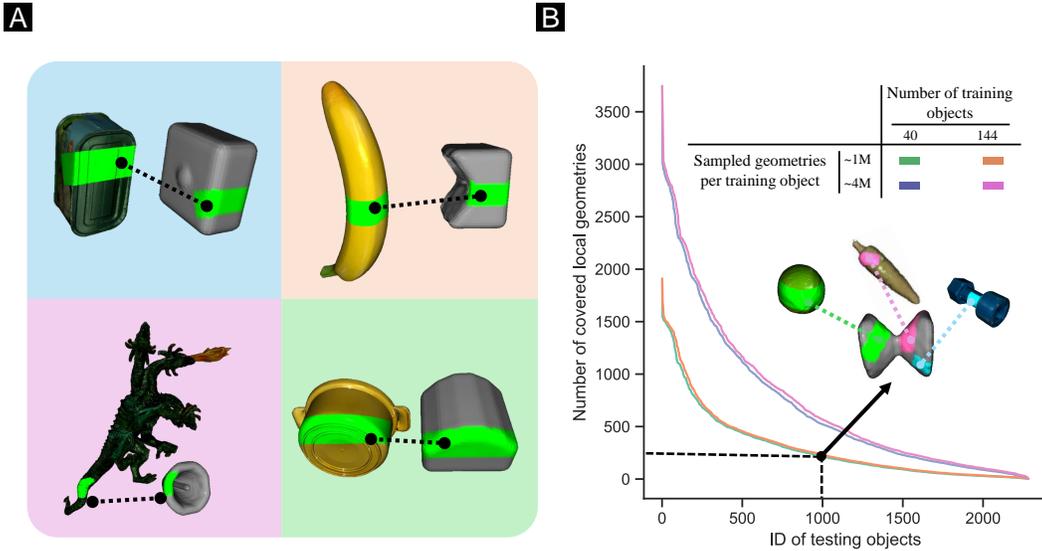


Figure 5: **Geometry coverage analysis.** (A): The colored objects are our training objects and the gray objects are in the EGAD! object set. The surfaces highlighted in green and connected by a dotted line have similar local geometries. We see that although the training object and testing object have very different overall shape, we can find local geometries on them that are pretty similar. (B): The local geometry coverage curves on the testing set given different choices of scaling up the training set. The x -axis denotes the ID of each testing object, and the y -axis denotes the number of covered local geometries on each testing object. An example is given where the 1000-th testing object has around 250 covered local geometries. We only draw 3 of them for illustration.

Based on this analysis, we prioritize scaling up the label density of CGRs on each training object, rather than increasing the number of training objects, when constructing our dataset. By training on over a billion CGRs, our model has learned to map local geometries to grasp representations effectively, thereby enhancing its ability to generalize to novel objects.

Now, we turn to the grasp decision model and discuss why it can learn grasp success from just hundreds of trial-and-error attempts. Here, we highlight a few possible reasons.

First, the representation captures all the relevant information about force closure that can be extracted from vision. For a point-to-plane contact problem, the force-closure condition must satisfy the following criteria Dai et al. (2018):

$$\begin{aligned}
 Gf &= 0, \\
 GG^T &> \epsilon I_{6 \times 6}, \\
 f_i^T n_i &> \frac{1}{\sqrt{\mu^2 + 1}} |f_i|,
 \end{aligned} \tag{1}$$

where f is the vector of contact forces acting at each contact point, G is the grasp matrix determined by the positions of the contact points, and n_i represents the surface normal at the i -th contact point. The latter two parameters are the only aspects that a vision model can estimate, and are generated by our representation model. The grasp decision model needs only to learn whether the forces f exerted by the gripper for different grasp types can satisfy Equation equation 1, given a friction coefficient μ . Although the friction coefficient is unknown, the model tends to learn an average behavior from the training set.

Second, the representation is compact. Instead of dealing with high-dimensional data like images or point clouds, we reduce the input to a 1D vector that represents the shape. This compactness simplifies the mapping from input to grasp quality, making it easier for the grasp decision model to learn.

4 MATERIALS AND METHOD

4.1 FORMULATION

A multi-finger grasp pose g is formally defined as:

$$g = [\mathbf{R} \ \mathbf{t} \ \mathbf{q}], \quad (2)$$

where $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ represents the robotic hand’s rotation, $\mathbf{t} \in \mathbb{R}^{3 \times 1}$ denotes the hand’s translation, and $\mathbf{q} \in \mathbb{R}^{n \times 1}$ characterizes the joint configuration of a n -DoF multi-finger hand. The goal of the grasp pose detection problem is to predict a set of grasp poses from a scene perception. Conventionally, data-driven methods have employed a single network $f(\cdot)$ to map the partial point cloud of the scene $\mathcal{P} \in \mathbb{R}^{k \times 3}$ to a set of candidate poses, $\mathbf{G} = \{g_i\}_{i=1}^{|\mathbf{G}|}$. In contrast, our approach decouples the mapping into two distinct steps: a state embedding step and a grasp decision step.

In the state embedding step, we extract a collection of contact-centric grasp representations from the partial point cloud \mathcal{P} . This is achieved by using a hand-agnostic representation model $\Phi(\cdot)$, which generates the scene representation \mathcal{R} :

$$\mathcal{R} = \Phi(\mathcal{P}), \quad (3)$$

where $\mathcal{R} = \{r_j\}_{j=1}^{|\mathcal{R}|}$ is a set of contact-centric grasp representations.

The grasp decision step consists of two distinct procedures: a mapping process that converts contact-centric grasp representations into a set of candidate grasp poses (referred to as grasp candidates) and a quality estimation process for each candidate. For each grasp representation r_j , we generate a set of grasp candidates based on the specific robotic hand. A hand-dependent mapping function $\mathcal{K}(\cdot)$ takes a grasp representation r_j and a hand specification h as input, and output \mathbf{G}_j :

$$\mathbf{G}_j = \mathcal{K}(r_j, h), \quad (4)$$

where $\mathbf{G}_j = \{g_j^{(i)}\}_{i=1}^{|\mathbf{G}_j|}$ denotes the set of grasp candidates for each r_j .

To estimate the quality of a grasp, we use a hand-dependent grasp decision model $\Psi(\cdot)$, which predicts the probability of success β given a grasp representation r_j , a grasp pose $g_j^{(i)}$, and a hand specification h :

$$\beta = \Psi(r_j, g_j^{(i)}, h). \quad (5)$$

Objective: Our goal is to find a set of grasp poses \mathbf{G}^* that maximizes the grasp success rate given a desired number of grasp poses K :

$$\mathbf{G}^* = \arg \max_{\mathbf{G} \subset \bigcup_j \mathbf{G}_j, |\mathbf{G}|=K} \mathbb{E}_{r_j \in \mathcal{R}, g_j^{(i)} \in \mathcal{K}(r_j, h) \cap \mathbf{G}} [\Psi(r_j, g_j^{(i)}, h)]. \quad (6)$$

4.2 CONTACT-CENTRIC GRASP REPRESENTATION

We initiate our approach with the development of a contact-centric grasp representation. Initially, consider a 2D object, we can represent it as a set comprising surface points and their corresponding normals:

$$r_{2d} = \{(p_i, n_i) \mid i = 1, 2, \dots, N\}. \quad (7)$$

In this representation, p_i denotes the position of a surface point, and n_i represents the normal vector associated with that surface point. For clarity, the object’s surface is discretized into N bins.

For the task of grasp pose detection, it is common to represent the object shape in a local coordinate frame ten Pas et al. (2017); Mousavian et al. (2019), as the classification of grasp quality depends primarily on the geometry within a localized area. This step, referred to as canonicalization, equips the representation with SE(3)-equivalent property and makes subsequent learning easier. For the 2D example, when we employ a polar coordinate system and sample the pole coordinate \mathbf{t}_{2d} and the polar axis \mathbf{R}_{2d} , the discrete object shape representation is refactored accordingly. In this system, a

surface point p_i is represented by an angle α_i from the polar axis and a distance d_i from the pole. Additionally, the surface normal is encoded as the angle between the normal n_i and α_i :

$$r_{2d} = \left\{ (\alpha_i, d_i, \theta_i) \mid i = 1, 2, \dots, N, p'_i = \mathbf{R}_{2d}(p_i - \mathbf{t}_{2d}), n'_i = \mathbf{R}_{2d}n_i, \right. \\ \alpha_i = \frac{p'_i}{\|p'_i\|}, d_i = \|p'_i\|, \\ \left. \theta_i = \arccos\left(\frac{\alpha_i \cdot n'_i}{\|\alpha_i\| \|n'_i\|}\right); \mathbf{R}_{2d}, \mathbf{t}_{2d} \right\}. \quad (8)$$

A benefit of adopting a polar coordinate system is that the in-plane rotation angles $\{\alpha_i\}$ can be uniformly sampled across the polar angle range, resulting in constant values for $\{\alpha_i\}$ across different representations. Therefore, we move $\{\alpha_i\}$ to the right side of the set notation to make the representation more compact. Since the values of d_i and θ_i depend on α_i , we rewrite them as d_{α_i} and θ_{α_i} :

$$r_{2d} = \left\{ (d_{\alpha_i}, \theta_{\alpha_i}) \mid \alpha_i = 0, \frac{2\pi}{N}, \dots, \frac{2\pi(N-1)}{N}; \mathbf{R}_{2d}, \mathbf{t}_{2d} \right\}. \quad (9)$$

Extending this representation to a real-world 3D object and a 3D coordinate system with rotation \mathbf{R}_{3d} and translation \mathbf{t}_{3d} involves decoupling the object’s geometry along a chosen axis and composing multiple 2D representations. By selecting a specific axis in the 3D coordinate system (*e.g.*, the z -axis), we discretize the object along this axis into M sections. Each section corresponds to a cross-sectional slice of the object at a particular coordinate along the axis.

Within each cross-sectional slice, the same polar coordinate system is employed as in the 2D case. We apply the same angular sampling and the local geometry is represented in terms of distance and normal angle at each sampled angle α_i . The 3D representation is then formulated as:

$$r_{3d} = \left\{ (d_{\alpha_i}, \theta_{\alpha_i})_j \mid \alpha_i = 0, \frac{2\pi}{N}, \dots, \frac{2\pi(N-1)}{N}, j = 1, 2, \dots, M; \mathbf{R}_{3d}, \mathbf{t}_{3d} \right\} \quad (10)$$

In the following sections, we use r as a shorthand for r_{3d} . In Figure 1B, we illustrate the process of representing a 3D geometry in the contact-centric representation format within a robotic hand’s local coordinate frame.

4.3 ROBOTIC HANDS AND GRASP TYPES

In our experiments, we utilize three distinct robotic hands:

- DH-3: A three-finger robotic hand comprises 4 degrees of freedom and 2 motors, operating in an underactuated manner.
- Allegro: A four-finger robotic hand comprises 16 degrees of freedom and 16 motors, designed for full actuation.
- Inspire: A five-finger robotic hand equipped with 12 degrees of freedom and 6 motors, operating in an underactuated manner.

These robotic hands represent a variety of applications, including industrial tasks, dexterous manipulation, and underactuated prosthetic hand functionalities.

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

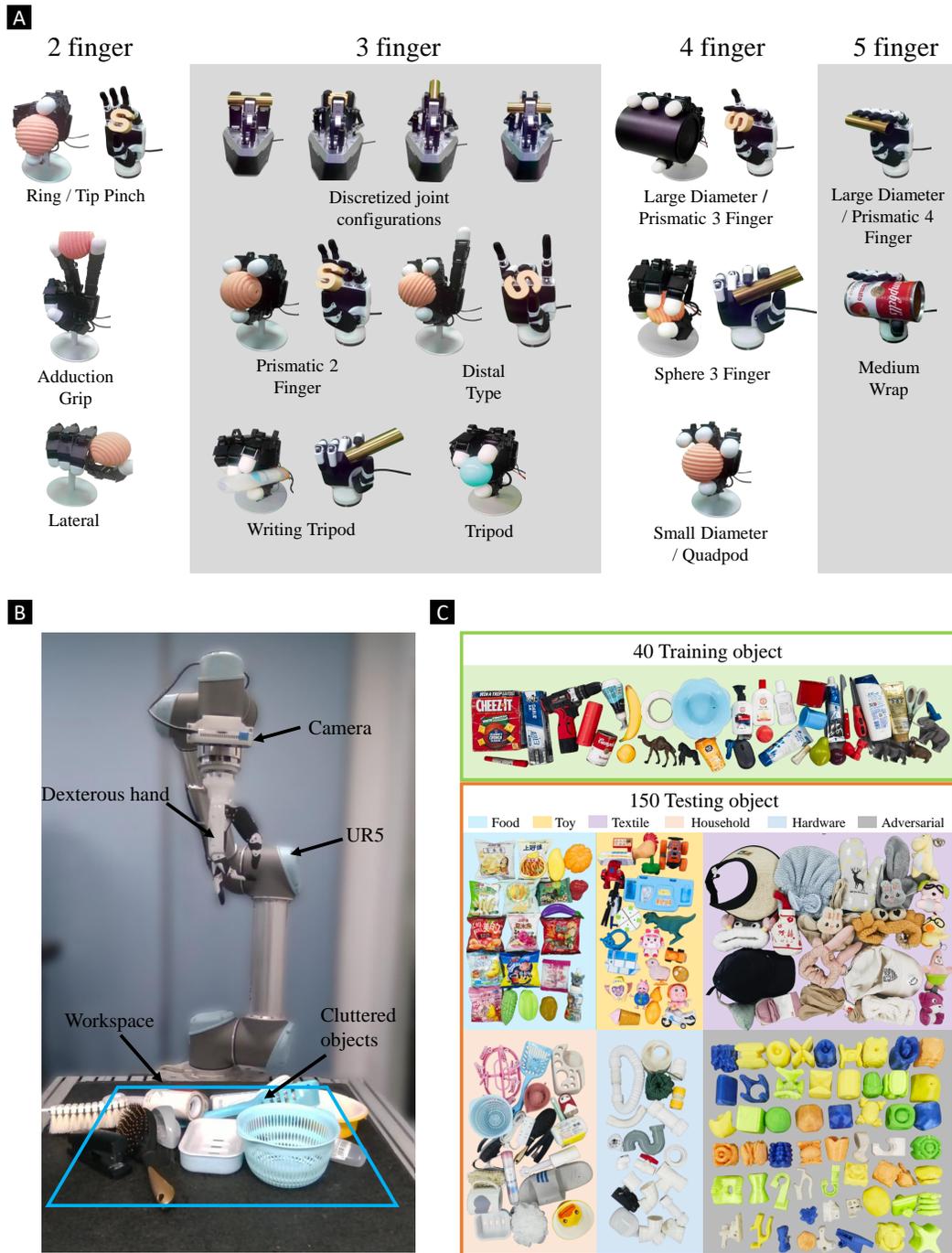


Figure 6: **Experimental setup.** (A): Illustration of the predefined grasp types for three robotic hands. The types are categorized by the number of fingers involved in the grasping procedure. Some types can be categorized into multiple taxonomies defined in previous literature Feix et al. (2015) when the grasping depths differ. (B): Platform setting of our dexterous grasping experiments. (C): Illustration of our 40 training objects and 150 testing objects. The testing objects are much more diverse than the training objects, including deformable and adversarial objects not presented in the training set.

756 One challenge with dexterous grasping is the complexity introduced by the high degrees of freedom
757 in these robotic hands, which creates a vast joint configuration space. However, when humans grasp
758 objects, we typically rely on only a small subset of these configurations, which can be categorized
759 into specific taxonomies Cutkosky et al. (1989); Feix et al. (2015). To address this complexity and
760 make grasp pose detection more manageable, we discretize the continuous joint configurations of
761 the multi-fingered hands into a finite set of predefined grasp types. This is represented as $\mathbf{q} \in$
762 $\{\mathbf{q}_1, \dots, \mathbf{q}_c\}$, where c denotes the total number of grasp types specific to each hand.

763 For the three-finger hand, we discretize the entire joint space into several bins, while for the four-
764 and five-finger hands, we select grasp types from the human grasp taxonomy that can be executed
765 by these dexterous robotic hands. The predefined grasp types are illustrated in Figure 6A. While
766 this approach simplifies the grasp pose detection process, it still provides sufficient flexibility for
767 subsequent manipulation tasks.

768 It is important to note that these grasp types serve as anchor poses prior to contact. Once the hand
769 reaches its target position, it undergoes a closure process, where the fingers progressively move
770 toward each other until the forces exerted on the finger joints reach predefined limits.

772 4.4 OVERVIEW OF EXPERIMENTS

774 To assess the performance of our multi-finger grasping model, we established a real-world exper-
775 imental platform. Our hardware setup includes a UR5 robotic arm and an Intel RealSense D415
776 camera, positioned at the robot’s end-effector. The initial camera pose is vertical to the table and is
777 approximately 60 cm above it. Figure 6B illustrates the setup of our robotic platform.

778 To assess the multi-finger grasp performance thoroughly, we construct a comprehensive real-world
779 test set, featuring objects commonly encountered in everyday life. These objects encompass diverse
780 shapes, materials, and textures and are categorized into hardware, food, textile, household, toy, and
781 adversarial items. The test set comprises nearly 150 objects ranging in size from $2.5 \times 2.5 \times 2.5$
782 cm^3 to $8 \times 8 \times 5 \text{cm}^3$.

783 During real-world testing, objects from each category are randomly placed on a table in a cluttered
784 way, and the robots attempt to grasp all the objects and clear the table. This process is repeated
785 twice for accuracy. We also establish a baseline that aligns the principal closing axis of the grasp
786 types with antipodal grasp poses, followed by collision detection, to compare with our proposed
787 method. The success rate is determined by dividing the number of successful grasp attempts by the
788 total number of grasp attempts.

789 Ultimately, our grasp system is successfully evaluated on three different robotic hands, where the
790 whole system is trained on a limited dataset comprising merely 40 objects and hundreds of grasp
791 attempts, and tested on a broader spectrum of 150 previously unseen objects. Notably, it represents
792 a pioneering achievement in the literature where a grasping algorithm is evaluated on a significantly
793 larger set of objects than those included in its training dataset. The final training and testing objects
794 are illustrated in Figure 6C for reference.

796 4.5 BASELINE METHOD

798 Here we introduce our baseline method of multi-finger grasping. Currently, our community can
799 achieve human-level robotic grasping with a parallel-jaw gripper Fang et al. (2023b). An intuitive
800 approach for multi-finger grasping is to mimic the behavior of parallel grasping. Thus, we propose
801 a baseline method that discovers the principal closing axis of a robotic hand and aligns it with a
802 parallel grasp pose. First, for each grasp type of a robotic hand, we manually designate its principal
803 closing axis, which is the primary direction along which the fingers converge when the hand closes
804 to grasp an object. Then, given a parallel grasp pose and a grasp type of a hand, we can align the
805 multi-finger hand’s principal closing axis to the parallel grasp pose. Previous literature Fan et al.
806 (2019; 2018) also explored similar ways to initialize a multi-finger grasp. In Figure 9 we illustrate
807 the alignment example.

808 When grasping with a selected robotic hand, we first generate multiple high-score antipodal grasp
809 poses for a single-view point cloud using the AnyGrasp library Fang et al. (2023b). Then, for each
antipodal grasp pose, we align the robotic hand configured in all grasp types with the antipodal pose.

810 It means that for each antipodal grasp pose, we would have multiple multi-finger grasp candidates
 811 with different types. For all of the multi-finger grasp candidates across the scene, we run a collision
 812 detection based on the partial-view point cloud and robotic hand model. We select the grasp type
 813 assigned with the highest antipodal grasp score for the remaining grasp candidates without collision.
 814 If multiple grasp candidates have the same grasp score, we randomly select one as the final grasp
 815 pose.

816 4.6 ALGORITHM DETAILS

817 Next, we introduce the details of our algorithm, which consists of three steps: learning the represen-
 818 tation model, mapping from representation to grasp pose, and learning the grasp decision model.
 819

820 4.6.1 REPRESENTATION MODEL

821 Our representation model takes a partial-view point cloud as input and generates the contact-centric
 822 grasp representation r for different rotation \mathbf{R}_{3d} and translation \mathbf{t}_{3d} across the scene. It may seem
 823 initially challenging to establish the representation model, given that this representation demands
 824 full surface information, and the r needs to be predicted for SE(3) space across the scene. How-
 825 ever, recent advancements in grasp pose detection algorithms have successfully learned the mapping
 826 from partial-view point clouds to antipodal grasp poses across the scene, unveiling the feasibility
 827 of learning the mapping from partial-view point clouds to the proposed intermediate representa-
 828 tion. Specifically, prior works, such as grasnet-baseline Fang et al. (2020b) and GSNet Wang et al.
 829 (2021), have predicted the gripper opening widths and antipodal scores for discretized rotation \mathbf{R}_{3d}
 830 and translation \mathbf{t}_{3d} across the scene:
 831

$$832 s = \left\{ (w_{\alpha_i}, \mu_{\alpha_i})_j \left| \alpha_i = 0, \frac{2\pi}{N}, \dots, \pi - \frac{2\pi}{N}, j = 1, 2, \dots, M; \mathbf{R}_{3d}, \mathbf{t}_{3d} \right. \right\}, \quad (11)$$

833 where $w_{\alpha_i} = 2 \times \max(d_{\alpha_i}, d_{\alpha_i+\pi})$ and $\mu_{\alpha_i} = \max(\tan(\theta_{\alpha_i}), \tan(\theta_{\alpha_i+\pi}))$ are the gripper opening
 834 width and antipodal grasp quality metric defined in Fang et al. (2020b). This representation shares a
 835 structural resemblance with our contact-centric grasp representation in Equation equation 10. Thus,
 836 we opt to build our representation model upon the GSNet Wang et al. (2021) architecture.
 837

838 When predicting the representation, it is intractable to account for every possible \mathbf{R}_{3d} and \mathbf{t}_{3d} in
 839 continuous space. Previous work Fang et al. (2020b); Wang et al. (2021) addressed this by selecting
 840 orientation \mathbf{R}_{3d} from 300 discretized directions, voxelizing the scene, and selecting only \mathbf{t}_{3d} that lies
 841 on object surfaces. However, the total number of resulting combinations still remains quite large.
 842 In Wang et al. (2021), a metric called “graspness” is proposed as a heuristic to bias sampling towards
 843 \mathbf{t}_{3d} and \mathbf{R}_{3d} values that have a higher probability of generating successful grasp poses. This metric
 844 includes two components: “point-wise graspness” and “view-wise graspness.” Point-wise graspness
 845 is calculated by counting the ratio of high-score antipodal grasp poses among all poses at a given \mathbf{t}_{3d} ,
 846 while view-wise graspness counts this ratio among all grasp poses at a given \mathbf{R}_{3d} for a sampled \mathbf{t}_{3d} .
 847 These two scores are learned jointly within the grasp pose detection network and guide sampling
 848 during inference.
 849

850 In this work, since we aim to train a hand-agnostic representation model, we define a new “grasp-
 851 ness” score for each r to indicate its suitability for subsequent grasping across different robotic
 852 hands. Intuitively, for a point $(\alpha_i, d_i, \theta_i)$, a robotic hand achieves better contact when θ_i is small,
 853 meaning the surface normal n_i is opposite to the contact direction (assuming the robotic finger ap-
 854 proaches towards the polar pole of the local coordinate frame). Additionally, geometries with many
 855 high-score antipodal grasp poses tend to be easier for dexterous hands to grasp. Thus, we define
 856 “graspness” in this paper as the sum of θ_i values below a threshold and the number of antipodal
 857 grasp poses in r . This definition helps the model reduce candidates for representation prediction
 858 within a scene without significantly affecting accuracy.

859 Similar to GSNet, our representation model consists of three cascaded modules. Firstly, a
 860 Minkowski Engine Choy et al. (2019) backbone takes the single-view point cloud as input, en-
 861 codes their geometric features, and outputs a computed feature vector for each input point. Then a
 862 multi-layer-perception (MLP) takes the features of each point and generates a point-wise graspness
 863 heatmap. We sample 1024 seed points with high graspness, and forward these points to another
 MLP block. It outputs the view-wise graspness scores for 300 approach directions towards each

864 seed point respectively. We then select the direction with the highest graspness score for each point,
 865 group the features with cylinder grouping Fang et al. (2020a) along that direction and forward the
 866 grouped features for each point through a final MLP block. This final layer outputs r for $N = 48$
 867 in-plane rotations and $M = 5$ grasp depths, which are 0.005m, 0.01m, 0.02m, 0.03m and 0.04m
 868 respectively.

869 4.6.2 MAPPING FROM REPRESENTATION TO MULTI-FINGER GRASP CANDIDATES

870 After we obtain the representation r at different positions across the scene, we link them with differ-
 871 ent grasp types of a robotic hand to generate multi-finger grasp candidates. In theory, since we have
 872 predicted the contact information, we can already generate suitable multi-finger grasp candidates
 873 through optimization Miller & Allen (2004); Liu et al. (2021). However, for simplicity, we follow
 874 the same technique adopted in the baseline method to generate multi-finger grasp candidates. Such
 875 a design also facilitates fair comparison with the baseline method and shows how our grasp decision
 876 network improves the grasping ability.
 877

878 For each predicted CGR with the form of Equation equation 10, we calculate the corresponding
 879 antipodal grasp representation defined in Equation equation 11. Then the CGRs with top-500 an-
 880 tipodal grasp scores are selected. These representations are associated with different multi-finger
 881 grasp candidates following the same procedure of the baseline method. After this process, we query
 882 the orientation \mathbf{R}_g and translation \mathbf{t}_g of the multi-finger grasp candidates associated with the CGRs
 883 (more details in supplementary material). Together with the associated grasp types, we map the
 884 CGRs to multi-finger grasp candidates.
 885

886 4.6.3 LEARNING MULTI-FINGER GRASPING

887 For each sampled grasp candidate g_i , we learn a mapping from its corresponding CGR r_i to grasp
 888 success probability. This mapping is approximated through the grasp decision model, using training
 889 data collected via trial and error:

$$890 \alpha = \Psi(r_i, g_i, h).$$

891 In practice, we train different decision models for different robotic hands, denoted as $\Psi_h(r_i, g_i)$.
 892 Since the grasp types of each hand are discretized, we further decompose the classification of differ-
 893 ent grasp types into different sub-models:
 894

$$895 \Psi_h(r_i, g_i) = \sum_{\mathbf{q}} \mathbb{I}(g_i, \mathbf{q}) \Psi_{h,\mathbf{q}}(r_i, \mathbf{R}_g, \mathbf{t}_g),$$

896 where $\mathbb{I}(g_i, \mathbf{q})$ is an indicator function that is 1 when g_i matches the grasp type \mathbf{q} and 0 otherwise.
 897 Since \mathbf{R}_g and \mathbf{t}_g are functions of r_i , we can simplify the input to the sub-models by removing \mathbf{R}_g
 898 and \mathbf{t}_g . Thus, $\Psi_{h,\mathbf{q}}(r_i, \mathbf{R}_g, \mathbf{t}_g)$ can be reformulated as $\Psi_{h,\mathbf{q}}(r_i)$, where the computation of \mathbf{R}_g and
 899 \mathbf{t}_g is implicit in the model. We empirically found that using different sub-models for different grasp
 900 types gives better performance. The input to the model consists of the CGR of a selected grasp
 901 candidate. The model’s output is a score of whether the selected grasp would be successful. Details
 902 of the model is given in supplementary material. For simplicity, we regard the combination of all
 903 sub-models for each robotic hand as a single model and still refers to it as the grasp decision model.
 904
 905

906 4.6.4 DETECTION POST-PROCESSING

907 Following the grasp decision model’s output, we select grasp poses with high-quality scores, typi-
 908 cally exceeding 0.9. Collision detection is then performed by voxelizing the pre-shaped multi-finger
 909 hand and examining intersections between the hand voxels and the scene point cloud using the
 910 Open3D library. The final grasp pose is chosen from those grasp poses without collision with the
 911 scene, with the highest grasp quality score.
 912

913 4.7 TRAINING ENVIRONMENT

914 4.7.1 TRAINING OBJECT SET

915 Our experiments involve three different training object sets, the larger dataset \mathbb{L} contains 144 objects,
 916 the smaller one \mathbb{S} contains 40 objects, and the tiniest one \mathbb{T} contains 30 objects. \mathbb{L} is the training
 917

Algorithm 1 Multi-finger Grasping Data Collection for Robotic Hand h **Input:** the expected size K of the grasp dataset.**Output:** the collected grasp dataset \mathbf{G} .

```

918 1:  $\mathbf{G} \leftarrow \emptyset$ 
919 2: while  $|\mathbf{G}| < K$  do
920 3:   The robot moves to the ready pose
921 4:    $\mathcal{P} \leftarrow \text{camera.perception}$   $\triangleright$  capture RGBD images and transform into point cloud
922 5:    $\mathcal{R} \leftarrow \Phi(\mathcal{P})$   $\triangleright$  generate scene representation from the point cloud
923 6:   Sample a CGR  $r \in \mathcal{R}$  in the scene
924 7:   Sample a grasp type  $\mathbf{q} \in \{\mathbf{q}_1, \dots, \mathbf{q}_c\}$ 
925 8:    $[\mathbf{R}_g \ \mathbf{t}_g] \leftarrow \text{compute\_grasp\_pose}(r)$   $\triangleright$  Map CGR to the grasp pose
926 9:   if  $\text{collision\_detection}([\mathbf{R}_g \ \mathbf{t}_g \ \mathbf{q}], \mathcal{P}; h)$  then
927 10:      $\triangleright$  Check if the multi-finger grasp pose will collide with the scene point cloud
928 11:     continue
929 12:   end if
930 13:   The robotic hand executes the multi-finger grasp pose  $[\mathbf{R} \ \mathbf{t} \ \mathbf{q}]$ 
931 14:   Record the grasp result  $S$   $\triangleright$  Collect trial-and-error results
932 15:    $\mathbf{G} \leftarrow \mathbf{G} \cup \{(r, \mathbf{R}_g, \mathbf{t}_g, \mathbf{q}, S)\}$ 
933 16: end while
934 17: return the collected grasp dataset  $\mathbf{G}$ 

```

set collected in AnyGrasp Fang et al. (2023b). \mathbb{S} encompasses the 40 training objects featured in the original GraspNet-1Billion dataset, and \mathbb{T} includes 30 randomly selected objects from \mathbb{L} . In Figure 7 we detail the three training object sets.

4.7.2 DATA ANNOTATION AND COLLECTION

To facilitate the training of our representation model, we re-annotate the GraspNet-1Billion dataset. The training set consists of 100 scenes made up of 40 objects. Each scene includes 256 RGBD images, each of which can be transformed into a single-view point cloud. Instead of the original antipodal grasp representation (illustrated in Equation equation 11), we annotate the contact-centric grasp representation as per Equation equation 10 for the 100 training scenes.

Our process begins by voxelizing the 3D mesh of each training object with a resolution of 0.005 m. We collect all points on the voxelized object surface, denoted as $\{\mathbf{t}_{3d}^{(i)}\}$, where i indexes each individual surface point. For each surface point $\mathbf{t}_{3d}^{(i)}$, we sample 300 approach directions $\{\mathbf{R}_{3d}^{(j)}\}$, where j indexes the sampled directions. We then compute the CGR r for each combination of $\mathbf{t}_{3d}^{(i)}$ and $\mathbf{R}_{3d}^{(j)}$. This computation relies on the complete mesh of the object. The computed CGRs are then projected from each training object to the training scenes based on the object’s 6D pose provided in the original dataset.

After generating the CGRs for each scene, we apply a simple post-processing step to verify grasp feasibility. For each CGR, we check whether a cylindrical region extending backward along the approach direction collides with the tabletop or other objects in the scene. If a collision is detected, we set the CGR to a zero vector, indicating it is not a viable grasp candidate. This post-processing step helps reduce the likelihood of robotic hand collisions within the scene.

To train the grasp decision model, we collect grasping data by trial and error. Previously, most of the grasp attempts related to multi-finger grasping were collected within a simulation environment. Nevertheless, significant gaps may arise due to the inherent differences between the simulation and real environments. Thus, in this paper, we directly collect grasping data in a real-world environment.

We provide an overview of the complete data collection pipeline, summarized in Algorithm 1. Initially, we randomly place objects on the table. We then run the representation model to generate dense contact-centric grasp representations for the scene. We sample a CGR and a grasp type of the robotic hand, and map the CGR to a multi-finger grasp pose. Collision detection is performed to ensure that the grasp pose does not collide with the scene. If no collision happens, we execute the

972 grasp process. During this process, we record whether the grasp is successful and store the necessary
973 information in the dataset.

975 4.7.3 TRAINING DETAILS

976
977 For the representation model, the input point clouds are down-sampled with a voxel size of 0.005m.
978 In practice, we set the parameters of the 3D representation N and M in Equation equation 10 to 48
979 and 5 respectively. The model is trained on the re-annotated GraspNet-1Billion dataset using one
980 Nvidia A100 GPU with Adam optimizer Kingma & Ba (2014) and an initial learning rate of 0.001.
981 The learning rate follows a descent strategy and we adopt “poly” policy with $power = 0.9$ for learn-
982 ing rate decay. The model is trained from scratch with a batch size of 4. For data augmentation,
983 we randomly flip the scene horizontally and randomly rotate the points by $\text{Uniform}[-30^\circ, 30^\circ]$
984 around the z -axis (in the camera coordinate frame). We also randomly translate the points by
985 $\text{Uniform}[-0.2\text{m}, 0.2\text{m}]$ in the x - or y -axis and $\text{Uniform}[-0.1\text{m}, 0.2\text{m}]$ in the z -axis.

986 For the grasp decision model, since we have a relatively limited amount of collected data, our model
987 is trained for only 20 epochs to avoid overfitting. We leverage the Adam optimizer Kingma & Ba
988 (2014). The learning rate follows a segmented descent strategy starting from 0.0001, and the batch
989 size \mathbf{Z} is set to 128 to optimize training efficiency. Since the network is quite small, we train the
990 model on a laptop with NVIDIA 1650 GPU.

991 4.8 EXPERIMENTAL PROCEDURE

992
993 In each experiment, we randomly distribute objects from different categories in the robot workspace.
994 During the grasping process, the partial-view point cloud captured by the camera is fed into our rep-
995 resentation model. When collecting training data, we follow the procedure in Algorithm 1. During
996 testing, we first choose 100 CGRs from the outcome of the representation model, which has the
997 top-100 antipodal grasp scores. These CGRs are mapped to multi-finger grasp candidates, and given
998 the number of predefined types for each robotic hand, the total number of multi-finger grasp candi-
999 dates varies (*e.g.*, we define 4 grasp types for the three-finger hand, thus it has 400 grasp candidates).
1000 These grasp candidates are fed into our grasp decision model. The grasp candidates with the top-200
1001 grasp quality scores then undergo collision detection post-processing. The grasp pose that passes
1002 collision detection and has the highest grasp score is selected as the final multi-finger grasp pose
1003 in the camera’s coordinate system. It is subsequently converted into the world coordinate system
1004 and sent to the UR5 robot through socket communication. The UR5’s embedded motion planner
1005 navigates it to the grasp pose, where we set a waypoint 10 cm backward from the final grasp along
1006 the approach direction to avoid collision during movement. Simultaneously, the robotic hand is
1007 configured to the selected grasp type. After the robot arm reaches the target pose, the robotic hand
1008 closes the fingers until the grasping force reaches a predefined limit. The robot then lifts the object
1009 and moves it to the top of the bin and drops the object. The experiment concludes with manually
1010 recording whether the robotic hand successfully move the object to target position.

1011 5 CONCLUSION

1012
1013 This paper proposed a novel methodology to address the challenges of dexterous grasping in clut-
1014 tered environments with minimal data requirements, demonstrating significant advancements in
1015 learning visually guided grasping across diverse robotic hands. By decoupling the grasping process
1016 into contact-centric representation learning and hand-specific grasp decision-making, the approach
1017 achieves robust real-world performance with minimal training—achieving a 75% to 95% success
1018 rate with just 40 training objects and hundreds of grasp attempts, and further improving with addi-
1019 tional training. The system’s adaptability to various robotic hands and its robustness with challeng-
1020 ing objects, including deformable and adversarial items, underscore its potential for broad applica-
1021 bility. This work establishes a foundation for efficient and scalable dexterous manipulation, paving
1022 the way for impactful use in domains such as humanoid robotics and prosthetics. Future directions
1023 for research entail expanding the scope of the contact-centric grasp representation model to include
1024 a wider array of tactile and sensory information, enabling a more comprehensive understanding of
1025 object manipulation. Tactile sensors encapsulate rich information concerning contact positions and
contact point normals, mirroring the fundamental attributes of our representation model.

REFERENCES

- 1026
1027
1028 Sergio Almecija, Salvador Moya-Sola, and David M Alba. Early origin for human-like precision
1029 grasping: a comparative study of pollical distal phalanges in fossil hominins. *PLoS One*, 5(7):
1030 e11727, 2010.
- 1031 Samarth Brahmhatt, Ankur Handa, James Hays, and Dieter Fox. Contactgrasp: Functional multi-
1032 finger grasp synthesis from contact. In *2019 IEEE/RSJ International Conference on Intelligent
1033 Robots and Systems (IROS)*, pp. 2386–2393. IEEE, 2019.
- 1034 Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski
1035 convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision
1036 and pattern recognition*, pp. 3075–3084, 2019.
- 1037
1038 Sammy Christen, Stefan Stevšić, and Otmar Hilliges. Guided deep reinforcement learning of control
1039 policies for dexterous human-robot interaction. In *2019 International Conference on Robotics and
1040 Automation (ICRA)*, pp. 2161–2167. IEEE, 2019.
- 1041 Enric Corona, Albert Pumarola, Guillem Alenya, Francesc Moreno-Noguer, and Grégory Rogez.
1042 Ganhand: Predicting human grasp affordances in multi-object scenes. In *Proceedings of the
1043 IEEE/CVF conference on computer vision and pattern recognition*, pp. 5031–5041, 2020.
- 1044 Mark R Cutkosky et al. On grasp choice, grasp models, and the design of hands for manufacturing
1045 tasks. *IEEE Transactions on robotics and automation*, 5(3):269–279, 1989.
- 1046
1047 Hongkai Dai, Anirudha Majumdar, and Russ Tedrake. Synthesis and optimization of force closure
1048 grasps via sequential semidefinite programming. *Robotics Research: Volume 1*, pp. 285–305,
1049 2018.
- 1050 Erik Domellöf, Marianne Barbu-Roth, Louise Rönnqvist, Anne-Yvonne Jacquet, and Jacqueline
1051 Fagard. Infant manual performance during reaching and grasping for objects moving in depth.
1052 *Frontiers in Psychology*, 6:1142, 2015.
- 1053
1054 Clemens Eppner, Arsalan Mousavian, and Dieter Fox. Acronym: A large-scale grasp dataset based
1055 on simulation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp.
1056 6222–6227. IEEE, 2021.
- 1057 Yongxiang Fan, Te Tang, Hsien-Chung Lin, and Masayoshi Tomizuka. Real-time grasp planning
1058 for multi-fingered hands by finger splitting. In *2018 IEEE/RSJ International Conference on In-
1059 telligent Robots and Systems (IROS)*, pp. 4045–4052. IEEE, 2018.
- 1060
1061 Yongxiang Fan, Xinghao Zhu, and Masayoshi Tomizuka. Optimization model for planning precision
1062 grasps with multi-fingered hands. In *2019 IEEE/RSJ International Conference on Intelligent
1063 Robots and Systems (IROS)*, pp. 1548–1554. IEEE, 2019.
- 1064 Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: A large-scale bench-
1065 mark for general object grasping. In *Proceedings of the IEEE/CVF conference on computer vision
1066 and pattern recognition*, pp. 11444–11453, 2020a.
- 1067
1068 Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: A large-scale bench-
1069 mark for general object grasping. In *Proceedings of the IEEE/CVF Conference on Computer
1070 Vision and Pattern Recognition*, 2020b.
- 1071 Hao-Shu Fang, Minghao Gou, Chenxi Wang, and Cewu Lu. Robust grasping across diverse sensor
1072 qualities: The graspnet-1billion dataset. *The International Journal of Robotics Research*, 42(12):
1073 1094–1103, 2023a.
- 1074
1075 Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhai Liu,
1076 Yichen Xie, and Cewu Lu. Anygrasp: Robust and efficient grasp perception in spatial and tem-
1077 poral domains. *IEEE Transactions on Robotics*, 2023b.
- 1078 Thomas Feix, Javier Romero, Heinz-Bodo Schmiedmayer, Aaron M Dollar, and Danica Kragic.
1079 The grasp taxonomy of human grasp types. *IEEE Transactions on human-machine systems*, 46
(1):66–77, 2015.

- 1080 Patrick Grady, Chengcheng Tang, Christopher D Twigg, Minh Vo, Samarth Brahmabhatt, and
1081 Charles C Kemp. Contactopt: Optimizing contact to improve grasps. In *Proceedings of the*
1082 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1471–1481, 2021.
- 1083 Abhishek Gupta, Clemens Eppner, Sergey Levine, and Pieter Abbeel. Learning dexterous manip-
1084 ulation for a soft robotic hand from human demonstrations. In *2016 IEEE/RSJ International*
1085 *Conference on Intelligent Robots and Systems (IROS)*, pp. 3786–3793. IEEE, 2016.
- 1086 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
1087 *arXiv:1412.6980*, 2014.
- 1088 Tracy L Kivell. Evidence in hand: recent discoveries and the early evolution of human manual
1089 manipulation. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 370(1682):
1090 20150105, 2015.
- 1091 Puhao Li, Tengyu Liu, Yuyang Li, Yiran Geng, Yixin Zhu, Yaodong Yang, and Siyuan Huang.
1092 Gendexgrasp: Generalizable dexterous grasping. In *2023 IEEE International Conference on*
1093 *Robotics and Automation (ICRA)*, pp. 8068–8074. IEEE, 2023.
- 1094 Min Liu, Zherong Pan, Kai Xu, Kanishka Ganguly, and Dinesh Manocha. Deep differentiable
1095 grasp planner for high-dof grippers. In *Proceedings of Robotics: Science and Systems (RSS)*,
1096 Cambridge, Massachusetts, 2020.
- 1097 Qingtao Liu, Yu Cui, Qi Ye, Zhengnan Sun, Haoming Li, Gaofeng Li, Lin Shao, and Jiming Chen.
1098 Dexreynet: Learning dexterous robotic grasping network with geometric and spatial hand-object
1099 representations. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems*
1100 *(IROS)*, pp. 3153–3160. IEEE, 2023.
- 1101 Tengyu Liu, Zeyu Liu, Ziyuan Jiao, Yixin Zhu, and Song-Chun Zhu. Synthesizing diverse and
1102 physically stable grasps with arbitrary hand structures using differentiable force closure estimator.
1103 *IEEE Robotics and Automation Letters*, 7(1):470–477, 2021.
- 1104 Jens Lundell, Francesco Verdoja, and Ville Kyrki. Ddgc: Generative deep dexterous grasping in
1105 clutter. *IEEE Robotics and Automation Letters*, 6(4):6899–6906, 2021.
- 1106 Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Apari-
1107 cio, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point
1108 clouds and analytic grasp metrics. In *Proceedings of Robotics: Science and Systems (RSS)*, Cam-
1109 bridge, Massachusetts, July 2017.
- 1110 Jeffrey Mahler, Matthew Matl, Vishal Satish, Michael Danielczuk, Bill DeRose, Stephen McKinley,
1111 and Ken Goldberg. Learning ambidextrous robot grasping policies. *Science Robotics*, 4(26),
1112 2019.
- 1113 Priyanka Mandikal and Kristen Grauman. Dexvip: Learning dexterous grasping with human hand
1114 pose priors from video. In *Conference on Robot Learning*, pp. 651–661. PMLR, 2022.
- 1115 Andrew T Miller and Peter K Allen. Graspit! a versatile simulator for robotic grasping. *IEEE*
1116 *Robotics & Automation Magazine*, 11(4):110–122, 2004.
- 1117 Douglas Morrison, Peter Corke, and Jürgen Leitner. Egad! an evolved grasping analysis dataset for
1118 diversity and reproducibility in robotic manipulation. *IEEE Robotics and Automation Letters*, 5
1119 (3):4368–4375, 2020.
- 1120 Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof graspnet: Variational grasp generation
1121 for object manipulation. In *International Conference on Computer Vision (ICCV)*, 2019.
- 1122 Karl M Newell, Deirdre M Scully, PV McDonald, and Renée Baillargeon. Task constraints and in-
1123 fant grip configurations. *Developmental Psychobiology: The Journal of the International Society*
1124 *for Developmental Psychobiology*, 22(8):817–831, 1989.
- 1125 Yuzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong
1126 Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. In *European*
1127 *Conference on Computer Vision*, pp. 570–587. Springer, 2022.

- 1134 Carlos Rosales, Lluís Ros, Josep M Porta, and Raúl Suárez. Synthesizing grasp configurations with
1135 specified contact regions. *The International Journal of Robotics Research*, 30(4):431–443, 2011.
1136
- 1137 Kenneth Shaw, Shikhar Bahl, Aravind Sivakumar, Aditya Kannan, and Deepak Pathak. Learning
1138 dexterity from human hand motion in internet videos. *The International Journal of Robotics
1139 Research*, 43(4):513–532, 2024.
- 1140 Matthew M. Skinner, Nicholas B. Stephens, Zewdi J. Tsegai, Alexandra C. Foote, N. Huynh
1141 Nguyen, Thomas Gross, Dieter H. Pahr, Jean-Jacques Hublin, and Tracy L. Kivell. Human-like
1142 hand use in australopithecus africanus. *Science*, 347(6220):395–399, 2015. doi: 10.1126/science.
1143 1261735.
- 1144 Andreas ten Pas, Marcus Gualtieri, Kate Saenko, and Robert Platt. Grasp pose detection in point
1145 clouds. *The International Journal of Robotics Research (IJRR)*, 36(13-14):1455–1473, 2017.
1146
- 1147 Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang. Unidex-
1148 grasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative
1149 generalist-specialist learning. *arXiv preprint arXiv:2304.00464*, 2023.
- 1150 Chenxi Wang, Hao-Shu Fang, Minghao Gou, Hongjie Fang, Jin Gao, and Cewu Lu. Graspness
1151 discovery in clutters for fast and accurate grasp detection. In *Proceedings of the IEEE/CVF
1152 International Conference on Computer Vision*, pp. 15964–15973, 2021.
- 1153 Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzhen Xu, Puhao Li, Tengyu Liu, and He Wang. Dex-
1154 graspnet: A large-scale robotic dexterous grasp dataset for general objects based on simulation.
1155 In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 11359–11366.
1156 IEEE, 2023.
- 1157 Wei Wei, Daheng Li, Peng Wang, Yiming Li, Wanyi Li, Yongkang Luo, and Jun Zhong. Dvvg: Deep
1158 variational grasp generation for dextrous manipulation. *IEEE Robotics and Automation Letters*, 7
1159 (2):1659–1666, 2022.
- 1160 Wei Wei, Peng Wang, Sizhe Wang, Yongkang Luo, Wanyi Li, Daheng Li, Yayu Huang, and Haonan
1161 Duan. Learning human-like functional grasping for multi-finger hands from few demonstrations.
1162 *IEEE Transactions on Robotics*, 2024.
- 1163 Yinzhen Xu, Weikang Wan, Jialiang Zhang, Haoran Liu, Zikang Shan, Hao Shen, Ruicheng Wang,
1164 Haoran Geng, Yijia Weng, Jiayi Chen, et al. Unidexgrasp: Universal robotic dexterous grasping
1165 via learning diverse proposal generation and goal-conditioned policy. In *Proceedings of the
1166 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4737–4746, 2023.
- 1167 Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing.
1168 *arXiv:1801.09847*, 2018.

1171 A SUPPLEMENTARY METHODS

1172 QUERY 6D GRASP POSE FOR CGR

1173
1174
1175
1176 When we map a CGR to grasp pose, we first calculate the antipodal grasp representation of the CGR.
1177 Given a CGR

$$1178 \quad r = \left\{ (d_{\alpha_i}, \theta_{\alpha_i})_j \mid \alpha_i = 0, \frac{2\pi}{N}, \dots, 2\pi - \frac{2\pi}{N}, j = 1, 2, \dots, M; \mathbf{R}_{3d}, \mathbf{t}_{3d} \right\}, \quad (12)$$

1180
1181 the antipodal grasp representation is calculated by

$$1182 \quad s = \left\{ (w_{\alpha_i}, \mu_{\alpha_i})_j \mid \alpha_i = 0, \frac{2\pi}{N}, \dots, \pi - \frac{2\pi}{N}, j = 1, 2, \dots, M; \mathbf{R}_{3d}, \mathbf{t}_{3d} \right\},$$

1183
1184
1185 where $w_{\alpha_i} = 2 \times \max(d_{\alpha_i}, d_{\alpha_i+\pi})$ and $\mu_{\alpha_i} = \max(\tan(\theta_{\alpha_i}), \tan(\theta_{\alpha_i+\pi}))$. After we obtained s ,
1186 we choose the α_i and j that has the maximum antipodal grasp score:

$$1187 \quad (\alpha_i^*, j^*) = \arg \max_{\alpha_i, j} (\mu_{\alpha_i})_j.$$

We add the rotation α_i^* and translation corresponds to the j^* -th section along approach direction to \mathbf{R}_{3d} and \mathbf{t}_{3d} :

$$\begin{aligned}\mathbf{R}_g &= \mathbf{R}_{3d} \cdot \mathbf{R}_z(\alpha_i^*), \\ \mathbf{t}_g &= \mathbf{t}_{3d} + \mathbf{d}(j^*) \cdot \mathbf{R}_{3d} \cdot \mathbf{R}_z(\alpha_i^*) \cdot \mathbf{z},\end{aligned}$$

where:

$$\mathbf{R}_z(\alpha_i^*) = \begin{bmatrix} \cos(\alpha_i^*) & -\sin(\alpha_i^*) & 0 \\ \sin(\alpha_i^*) & \cos(\alpha_i^*) & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

$\mathbf{d}(\cdot)$ is a function that maps the index j of the section to its actual depth along the approach direction (maps $\{1, 2, 3, 4, 5\}$ to $\{0.005\text{m}, 0.01\text{m}, 0.02\text{m}, 0.03\text{m}, 0.04\text{m}\}$), and:

$$\mathbf{z} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

The updated rotation \mathbf{R}_g and translation \mathbf{t}_g are the 6D grasp pose that corresponds to this CGR.

DETAILS OF GRASP DECISION MODEL

Each grasp decision sub-model is learned by a neural network. It takes a contact-centric grasp representation as input and outputs a score ranging from 0 to 1 to indicate whether the corresponding grasp candidate would be successful. The input size is $2 \times 5 \times 48 = 480$, which is composed of distances and normal angles on 5 sections along 48 in-plane rotations.

The network comprises seven fully connected layers with a skip connection for improving robustness. Each intermediate layer consists of a fully connected layer with 1024 neurons, a batch normalization layer, and a ReLU activation function. The output of the second intermediate layer is also forwarded to the fifth intermediate layer with a skip connection. Networks for different grasp types are trained separately. We employ a loss function defined as:

$$L = -\frac{1}{Z} \sum_{z=1}^Z y_z \log(p_z). \quad (13)$$

In this equation, L is the loss, y denotes the binary label of whether the real robot trial and error succeeded or not, and p represents the predicted grasp success probability by the network. Z denotes the batch size.

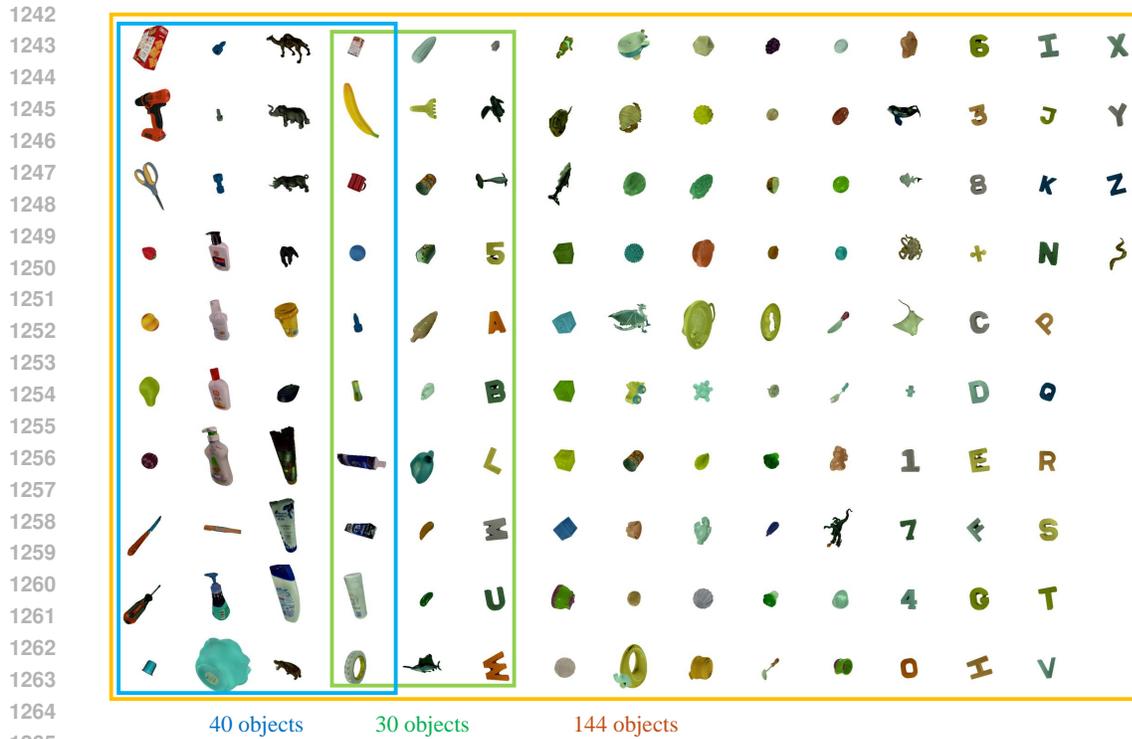
LOCAL GEOMETRY SAMPLING FOR GRASP COVERAGE ANALYSIS

The local geometries are cropped using 3D boxes defined by valid antipodal grasp poses. To obtain the grasp candidates, each object is voxel-downsampled to get grasp points in uniform distributions. V approach directions are sampled on the grasp point. A inplane rotation angles are sampled uniformly for each direction. On the training objects, we set $V=100$ and $A=12$ for the dense set, and $V=50$ and $A=6$ for the sparse set, respectively. In these two cases, the average numbers of local geometries for each training object are around 1M and 4M. For testing objects in the EGAD dataset, we set $V=100$ and $A=12$.

B SUPPLEMENTARY TEXT

TRAINING OBJECT COLLECTION

The 144 training objects are collected from supermarkets and grocery stores, which is extended from the 40 training objects collected in GraspNet-1Billion Fang et al. (2023a). The principle of choosing objects is that they have a roughly different shape or some local geometries from other objects, and they are chosen by authors heuristically. We provide the 3D scanned models of the objects to support reproducible research. Figure 7 shows an overview of the training object.



1267 **Figure 7: Training object set.** The set \mathbb{L} with 144 training objects are enclosed by the orange
1268 rectangle, the set \mathbb{S} with 40 training objects are enclosed by the blue rectangle, and the set \mathbb{T} with 30
1269 training objects are enclosed by the green rectangle. Their CAD models are available upon request.

1271 GRASP TYPES FOR DIFFERENT HANDS

1272 The index number for each grasp type of different robotic hands is given in Figure 8.

1275 PRINCIPAL CLOSING AXIS FOR DIFFERENT GRASP TYPES

1276 We illustrate the principal closing axis for different grasp types in Figure 9. The x -axis (in red) in
1277 the local coordinate frame is the approach direction and the y -axis (in green) is the principal closing
1278 axis of the hand. The two-finger gripper (in blue) is the corresponding antipodal grasp pose for each
1279 grasp type.
1280

1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

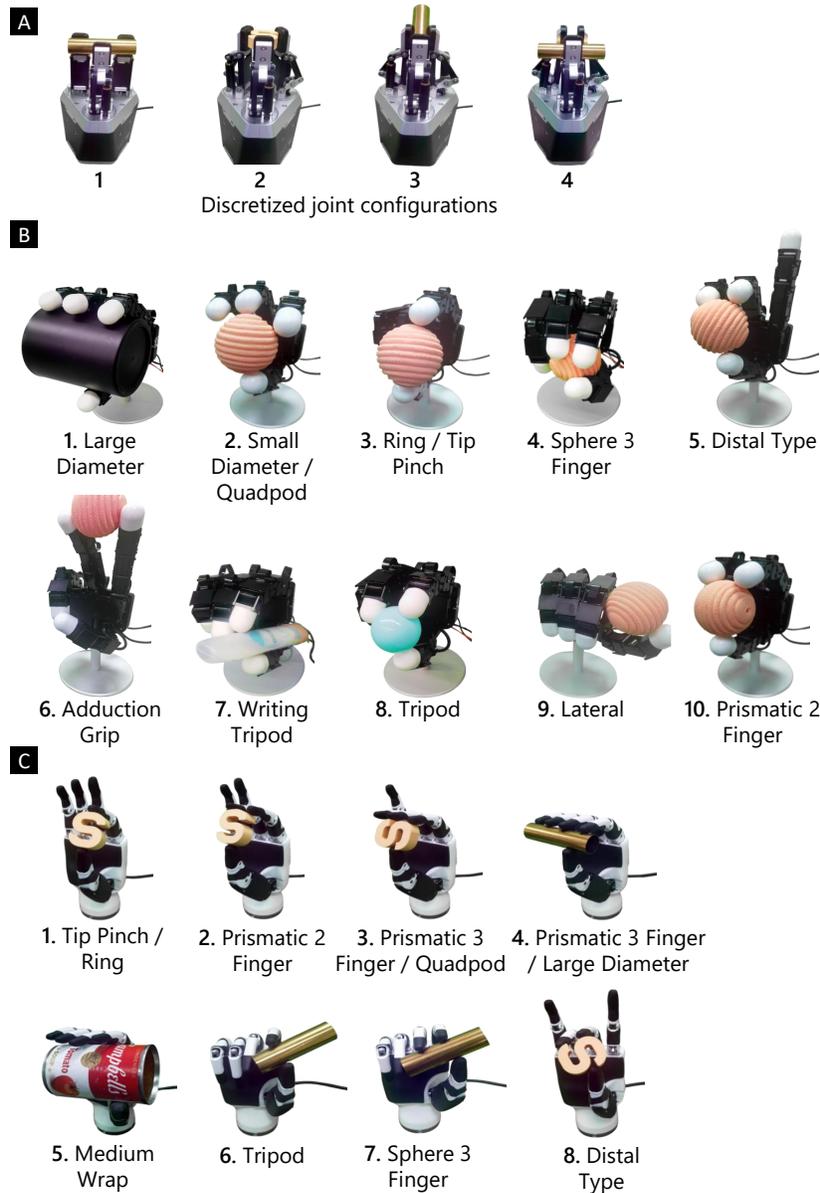


Figure 8: **Grasp type numbering.** (A), (B) and (C) give the index numbers of different grasp types for the three-finger, four-finger, and five-finger hands, respectively.

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

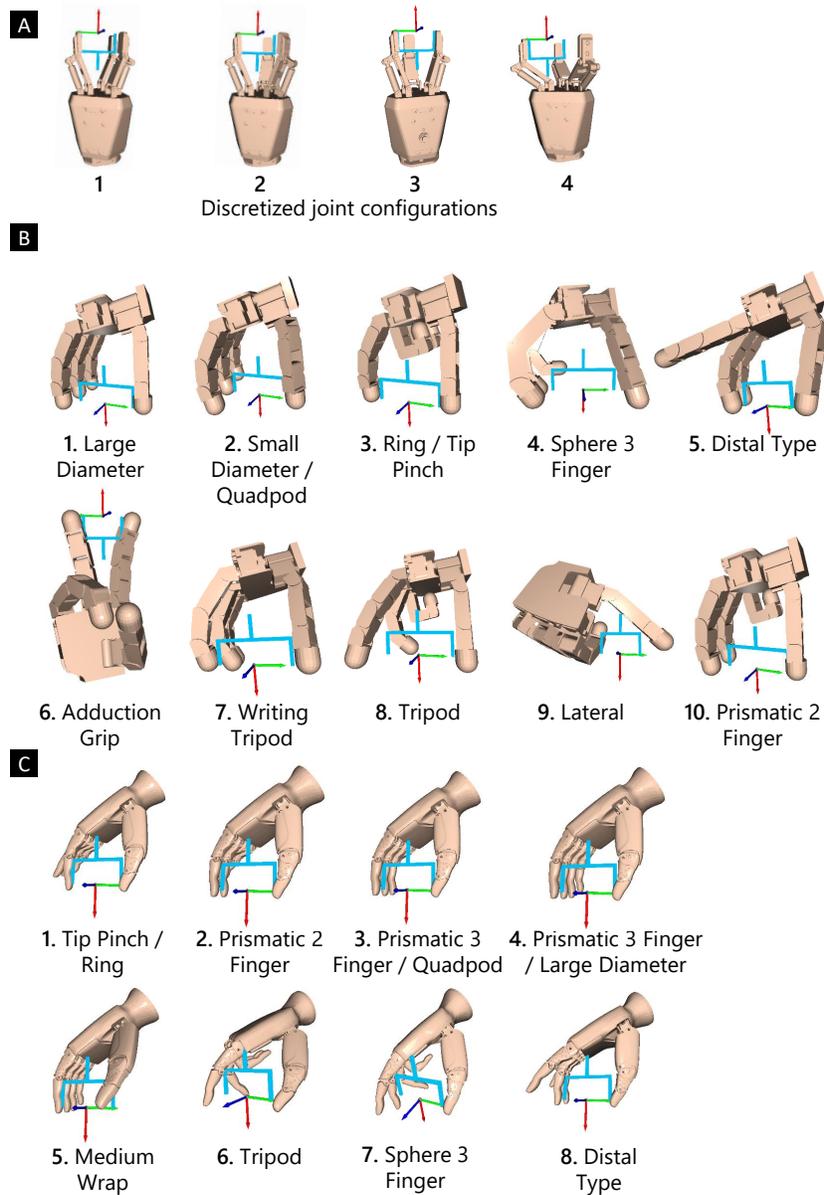


Figure 9: **Identifying principal closing axis.** We show the designated principal closing axis and corresponding antipodal grasp pose for each predefined grasp type. (A), (B) and (C) shows the results for the three hands we used.

C SUPPLEMENTARY MOVIES

We believe that presenting the complete process of our robotic grasping experiments can provide valuable insights into potential improvements for the grasping system. Additionally, it is essential to demonstrate the system’s robustness, which requires running it for an extended period. Therefore, we recorded the entire grasping process, retaining all original content without cuts, but with speed adjustments to keep the video at a reasonable length. The grasping process for each robotic hand lasts over 3 hours, with the total time across all three hands exceeding 15 hours. We applied a 20x speed-up for the collision detection phase and a 2x speed-up for the grasp execution phase. Even after these adjustments, the resulting videos still exceed 6 hours in length. Consequently, we have hosted the videos on YouTube, with the links provided below:

- Movie S1 - Grasping with 3-finger DH-3 hand on daily objects, after training on 144 objects:
<https://youtu.be/GGBesshyfxk>
- Movie S2 - Grasping with 4-finger Allegro hand on daily objects, after training on 144 objects:
https://youtu.be/HkrvWm_TTGo
- Movie S3 - Grasping with 5-finger Inspire hand on daily objects, after training on 144 objects:
<https://youtu.be/3Om7G8nMJPg>
- Movie S4 - Grasping with 3-finger DH-3 hand on adversarial objects, after training on 144 objects:
<https://youtu.be/GGBesshyfxk?t=1837>
- Movie S5 - Grasping with 4-finger Allegro hand on adversarial objects, after training on 144 objects:
<https://youtu.be/E7i3pqxA4RM>
- Movie S6 - Grasping with 5-finger Inspire hand on adversarial objects, after training on 144 objects:
<https://youtu.be/o6LQwRgu82s>
- Movie S7 - Grasping with 3-finger DH-3 hand on daily and adversarial objects, after training on 40 objects: <https://youtu.be/--5wIHfPoZs>
- Movie S8 - Grasping with 4-finger Allegro hand on daily objects, after training on 40 objects:
<https://youtu.be/uhaC8NORqm4>
- Movie S9 - Grasping with 4-finger Allegro hand on adversarial objects, after training on 40 objects:
<https://youtu.be/5pN6BYOH4xw>
- Movie S10 - Grasping with 5-finger Inspire hand on daily objects, after training on 40 objects:
<https://youtu.be/GQDLTVjXPQk>
- Movie S11 - Grasping with 5-finger Inspire hand on adversarial objects, after training on 40 objects:
<https://youtu.be/B7qc7qRw4ss>