

# $\mathcal{D}(\mathcal{R}, \mathcal{O})$ Grasp: A Unified Representation of Robot and Object Interaction for Cross-Embodiment Dexterous Grasping

Anonymous Author(s)

Affiliation

Address

email

1           **Abstract:** Dexterous grasping is a fundamental yet challenging skill in robotic  
2 manipulation, requiring precise interaction between robotic hands and objects. In  
3 this paper, we present  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  Grasp, a novel framework that models the interac-  
4 tion between the robotic hand in its grasping pose and the object, enabling broad  
5 generalization across various robot hands and object geometries. Our model takes  
6 the robot hand’s description and object point cloud as inputs and efficiently pre-  
7 dicted kinematically valid and stable grasps, demonstrating strong adaptability to  
8 diverse robot embodiments and object geometries. Extensive experiments con-  
9 ducted in both simulated and real-world environments validate the effectiveness  
10 of our approach, with significant improvements in success rate, grasp diversity,  
11 and inference speed across multiple robotic hands. Our method achieves an av-  
12 erage success rate of **87.53%** in simulation in less than one second, tested across  
13 three different dexterous robotic hands. In real-world experiments using the Leap-  
14 Hand, the method also demonstrates an average success rate of **89%**.  $\mathcal{D}(\mathcal{R}, \mathcal{O})$   
15 Grasp provides a robust solution for dexterous grasping in complex and varied en-  
16 vironments. The code, appendix, and videos are available on our project website  
17 at <https://drograsp.github.io/>.

18           **Keywords:** Dexterous Grasping, Robotic Manipulation

## 19   1 Introduction

20           Dexterous grasping is crucial in robotics as the first step in  
21 executing complex manipulation tasks. However, quickly  
22 obtaining a high-quality and diverse set of grasps re-  
23 mains challenging for dexterous robotic hands due to their  
24 high degrees of freedom and the complexities involved in  
25 achieving stable, precise grasps. Researchers have de-  
26 veloped several optimization-based methods to address this  
27 challenge [1, 2, 3, 4, 5]. Some of these methods, however,  
28 often focus on fingertip point contact, relying on com-  
29 plete object shape, and require significant computational  
30 time to optimize. As a result, data-driven grasp genera-  
31 tion methods have gained attention. These methods aim  
32 to solve the grasping problem using learning-based tech-  
33 niques. We can broadly categorize them into two types:  
34 those that utilize robot-centric representations, such as  
35 wrist poses and joint values [6, 7, 8], and those that rely  
36 on object-centric representations, such as contact points [9, 10, 11] or contact maps [12, 13, 14, 15].

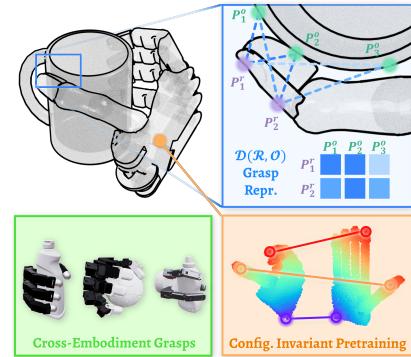


Figure 1: We propose our model that utilizes configuration-invariant pre-  
training, predicts  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation, and obtains grasps for cross-  
embodiment from point cloud input.

	Grasp Representation	Method Type	Cross Embodiment	Inference Speed	Sample Efficiency	Partial Object Point Cloud	Full-hand Contact (not only fingertips)	Optional Grasp Preference Interface
DFC [2]	Joint Values	Robot-centric	✓	XX	-	X	✓	X
UniDexGrasp++ [8]	Joint Values	Robot-centric	X	✓✓	X	✓	✓	X
UniGrasp [9]	Contact Point	Object-centric	✓	X	✓	X	X	X
GeoMatch [10]	Contact Point	Object-centric	✓	X	✓	X	✓	X
GenDexGrasp [12]	Contact Map	Object-centric	✓	X	✓	X	✓	X
ManiFM [13]	Contact Map	Object-centric	✓	X	✓	X	X	X
<b>DRO-Grasp (Ours)</b>	$\mathcal{D}(\mathcal{R}, \mathcal{O})$	Interaction-centric	✓	✓	✓	✓	✓	Contact Region Palm Orientation

Table 1: Dexterous grasp method comparison.

37 Robot-centric representations (e.g., joint values), as used in methods like UniDexGrasp++ [8], directly map observation to control commands for fast inference but suffer from low sample efficiency  
38 and poor generalization across different robot embodiments. The learned mappings are specific to  
39 the training data and do not quickly adapt to new robot designs or geometries. Object-centric rep-  
40 resentations (e.g., key points, contact points, affordances) effectively capture the geometry and con-  
41 tacts of objects, allowing for generalization across different shapes and robots, as demonstrated by  
42 methods like UniGrasp [9] and GenDexGrasp [12]. However, these methods are often less efficient  
43 as they typically require an additional optimization step—such as solving fingertip inverse kine-  
44 matics (IK) or fitting the predicted contact maps under penetration-free and joint limit constraints  
45 to translate the object-centric representation into actionable robot commands. This optimization  
46 process is time-consuming due to its complexity and nonconvexity [16, 17, 18].

47 To overcome the limitations of both paradigms, we propose  $\mathcal{D}(\mathcal{R}, \mathcal{O})$ , a unified representation that  
48 captures the relationship between the robotic hand’s grasp shape and the object.  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  encap-  
49 sulates both the articulated structure of the robot hand and the object’s geometry, enabling direct  
50 inference of kinematically valid and stable grasps that generalize across various shapes and robot  
51 embodiments.

52 Given the point clouds of both an open robotic hand and the object, our network architecture predicts  
53 the  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation. This matrix encodes the relative distances between the point clouds of  
54 the object and the robotic hand in the desired grasping pose[19, 20]. Using this representation, we  
55 apply a multilateration method [21] to estimate the robot’s point cloud at the predicted pose, allowing  
56 us to compute the 6D pose of each hand link in the world frame and ultimately determine the joint  
57 configurations. To encode robotic hands, we propose a configuration-invariant pretraining method  
58 that learns the inherent alignment between various hand configurations, promoting grasp generation  
59 performance and cross-embodiment generalization. We validate the effectiveness of our approach  
60 through extensive experiments in both simulation and real-world settings. Our model achieves an  
61 average success rate of 87.53% in simulation across three dexterous robotic hands and in real-robot  
62 experiments, demonstrating its robustness and versatility.

63 In conclusion, our primary contributions are as follows:

- 65 1. We introduce a novel representation,  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  for dexterous grasping tasks. This interaction-  
66 centric formulation transcends conventional robot-centric and object-centric paradigms, facil-  
67 itating robust generalization across diverse robotic hands and objects.
- 68 2. We propose a configuration-invariant pretraining approach with contrastive learning, establish-  
69 ing inherent alignment across varying configurations of robotic hands. This unified task can  
70 facilitate valid grasp generation and cross-embodiment feature alignment.
- 71 3. We perform extensive experiments in both simulation environments and real-world settings,  
72 validating the efficacy of our proposed representation and framework in grasping novel objects  
73 with multiple robotic hands.

## 74 2 Method

75 Given the object point cloud and the robot hand URDF file, our goal is to generate dexterous and  
76 diverse grasping poses that generalize across various objects and robot hands. Fig. 2 provides an  
77 overview of our proposed method.

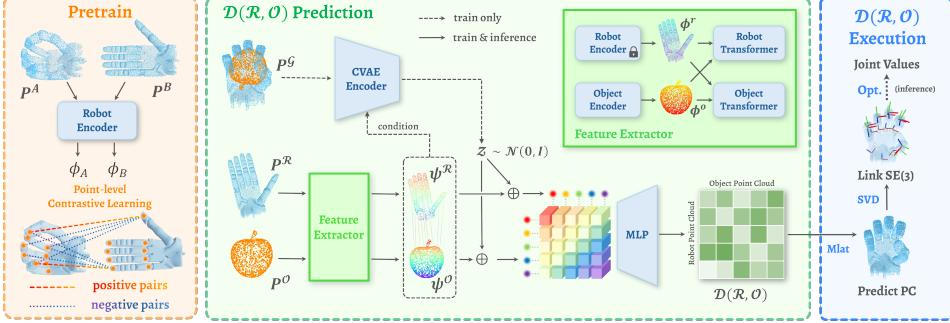


Figure 2: Overview of  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  framework: We first pretrain the robot encoder with the proposed configuration-invariant pretraining method. Then, we predict the  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation between the robot and object point cloud. Finally, we extract joint values from the  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation.

78 **Method Overview.** First, we design an encoder network to learn representations from the point  
 79 clouds of both the robot and the object. The robot encoder network is pretrained using our proposed  
 80 configuration-invariant pretraining method (Sec. 2.1), which facilitates the learning of efficient robot  
 81 embedding. Next, a CVAE model is used to predict the  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation, a point-to-point  
 82 distance matrix between the robotic hand at its grasp pose and the object, to implicitly present the  
 83 grasp pose (Sec. 2.2). From the  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation, we derive the 6D pose for each link,  
 84 which serves as the optimization target for determining the joint values. This optimization process  
 85 is notably straightforward and efficient (Sec. 2.3).

## 86 2.1 Configuration-Invariant Pretraining with Contrastive Learning

87 Learning dexterous grasping involves understanding the spatial relation-  
 88 ships between the robot hand and the object. The objective is to match  
 89 the robot hand in a specific configuration with the object. However, this  
 90 matching process is challenging because the local geometric features of  
 91 a point in the open-hand configuration may not align with those in the  
 92 grasp configuration due to huge variations during articulation.

93 To address this, we break the problem into two simpler components: (1)  
 94 self-articulation matching, which implicitly determines the joint values  
 95 for the grasp configuration, and (2) wrist pose estimation. As shown in  
 96 Fig. 3, leveraging configuration-invariant pretraining, we train the neu-  
 97 ral network to understand the self-articulation alignment across differ-  
 98 ent configurations, thereby facilitating the matching process between the  
 99 robot hand and the object.

100 Specifically, for each robot hand, we begin by uniformly sampling points  
 101 on the surface of each link at the canonical pose, storing the resulting  
 102 point clouds denoted as  $\{\mathbf{P}_{\ell_i}\}_{i=1}^{N_\ell}$ , where  $N_\ell$  is the number of links. We  
 103 define a point cloud forward kinematics model,  $\text{FK}(q, \{\mathbf{P}_{\ell_i}\}_{i=1}^{N_\ell})$  to map joint configura-  
 104 tions to point clouds at new poses. For example, given a close-hand  $q_A$  and an open-hand configura-  
 105 tion  $q_B$ , where the wrist pose is the same or nearly identical, we obtain two point clouds  $\mathbf{P}^A, \mathbf{P}^B \in \mathbb{R}^{N_R \times 3}$ ,  
 106 representing these two joint configurations. Here,  $N_R$  is the number of points in the robot point  
 107 cloud, set to 512 in practice.

108 These point clouds are passed through the encoder network (as described in Sec. 2.2) to produce  
 109 point-wise features  $\phi^A, \phi^B \in \mathbb{R}^{N_R \times D}$ , where  $D = 512$  is the feature dimension. The model  
 110 applies point-level contrastive learning, aligning embeddings of positive pairs—points with the same  
 111 index in both clouds—while separating negative pairs, weighted by the Euclidean distance in  $\mathbf{P}^B$ .  
 112 This process ensures that the features corresponding to the same positions on the robot hand remain

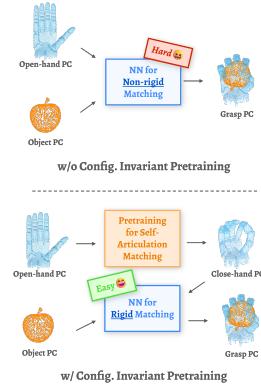


Figure 3: Motivation for configuration-invariant pretraining.

113 consistent across different joint configurations. We define the resulting contrastive loss as:

$$\mathcal{L}_p = -\frac{1}{N_\ell} \sum_i \log \left[ \frac{\exp(\langle \phi_i^A, \phi_i^B \rangle / \tau)}{\sum_j \omega_{ij} \exp(\langle \phi_i^A, \phi_j^B \rangle / \tau)} \right], \quad (1)$$

$$\omega_{ij} = \begin{cases} \frac{\tanh(\lambda \| p_i^B - p_j^B \|_2)}{\max(\tanh(\lambda \| p_i^B - p_j^B \|_2))}, & \text{if } i \neq j \\ 1, & \text{if } i = j \end{cases}, \quad (2)$$

114 where  $\langle \cdot, \cdot \rangle$  denotes the cosine similarity between two vectors,  $p_i^B$  represents the  $i$ -th point position  
115 in  $\mathbf{P}^B$ . For the hyperparameters, we set  $\tau = 0.1$  and  $\lambda = 10$  in practice. Note that the learned  
116 features are finger configuration-invariant but dependent on the wrist pose.

## 117 2.2 $\mathcal{D}(\mathcal{R}, \mathcal{O})$ Prediction

118 Denote an open-hand configuration as  $q_{init}$ , of which the wrist pose can be either user-specified  
119 or randomly generated. Let the robot point cloud under  $q_{init}$  be  $\mathbf{P}^R = \text{FK}(q_{init}, \{\mathbf{P}_{\ell_i}\}_{i=1}^{N_\ell}) \in$   
120  $\mathbb{R}^{N_R \times 3}$ , and the object point cloud be  $\mathbf{P}^O \in \mathbb{R}^{N_O \times 3}$ , where  $N_O$  represents the number of points in  
121 the object point cloud, also set to 512 in practice. The objective of our neural network is to predict  
122 the point-to-point distance matrix  $\mathcal{D}(\mathcal{R}, \mathcal{O}) \in \mathbb{R}^{N_R \times N_O}$ .

123 **Point Cloud Feature Extraction** We begin by extracting point cloud embeddings using two en-  
124 coders,  $f_{\theta_R}(\mathbf{P}^R)$  and  $f_{\theta_O}(\mathbf{P}^O)$ , which share the same architecture. Specifically, we use a modified  
125 DGCNN [22] to better capture local structures and integrate global information (see Appendix). The  
126 robot encoder is initialized with pretrained parameters, using the method described in Sec. 2.1, and  
127 remains frozen during training. These encoders extract point-wise features,  $\phi^R$  and  $\phi^O$  from the  
128 robot and object point clouds:

$$\phi^R = f_{\theta_R}(\mathbf{P}^R) \in \mathbb{R}^{N_R \times D}, \phi^O = f_{\theta_O}(\mathbf{P}^O) \in \mathbb{R}^{N_O \times D}. \quad (3)$$

129 To establish correspondences between the robot and object features, we apply two multi-head cross-  
130 attention transformers [23] (see Appendix),  $g_{\theta_R}(\phi^R, \phi^O)$  and  $g_{\theta_O}(\phi^O, \phi^R)$ . These transformers  
131 integrate the relationships between the two feature sets, embedding correspondence information.  
132 This process maps the robot and object features to two sets of correlated features,  $\psi^R$  and  $\psi^O$ :

$$\psi^R = g_{\theta_R}(\phi^R, \phi^O) + \phi^R \in \mathbb{R}^{N_R \times D}, \psi^O = g_{\theta_O}(\phi^O, \phi^R) + \phi^O \in \mathbb{R}^{N_O \times D}. \quad (4)$$

133 **CVAE-based  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  Prediction** To achieve cross-embodiment grasp diversity, we employ a  
134 Conditional Variational Autoencoder (CVAE) [24] network to capture variations across numerous  
135 combinations of hand, object, and grasp configurations. The CVAE encoder  $f_{\theta_g}$  takes the robot  
136 and object point clouds under the grasp pose  $\mathbf{P}^G \in \mathbb{R}^{(N_R+N_O) \times 3}$ , along with the learned features  
137  $(\psi^R, \psi^O)$ , resulting in an input shape of  $(N_R + N_O) \times (3 + D)$ . The encoder outputs the latent  
138 variable  $z \in \mathbb{R}^d$ , set as  $d = 64$  in practice. We concatenate  $z$  with extracted features  $\psi^R$  and  $\psi^O$ ,  
139 converting the feature to  $\hat{\psi}_i^R, \hat{\psi}_j^O \in \mathbb{R}^{N_O \times (D+d)}$ .

140 The same kernel function  $\mathcal{K}$  as Eisner et al. [25] is adopted, which possesses the properties of non-  
141 negativity and symmetry, to predict pair-wise distance  $r_{ij} = \mathcal{K}(\hat{\psi}_i^R, \hat{\psi}_j^O) \in \mathbb{R}^+$  under the grasp  
142 pose:

$$\mathcal{K}(\hat{\psi}_i^R, \hat{\psi}_j^O) = \sigma \left( \frac{1}{2} \mathcal{N}_\theta \left( \hat{\psi}_i^R, \hat{\psi}_j^O \right) + \frac{1}{2} \mathcal{N}_\theta \left( \hat{\psi}_j^O, \hat{\psi}_i^R \right) \right), \quad (5)$$

143 where  $\sigma$  denotes the softplus function, and  $\mathcal{N}_\theta$  is an MLP, which takes in the feature of  $\mathbb{R}^{N_O \times (2D+2d)}$   
144 and outputs a positive number (see Appendix). By calculating on all  $(\hat{\psi}_i^R, \hat{\psi}_j^O)$  pairs, we obtain the  
145 complete  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation:

$$\mathcal{D}(\mathcal{R}, \mathcal{O}) = \begin{bmatrix} \mathcal{K}(\hat{\psi}_1^R, \hat{\psi}_1^O) & \cdots & \mathcal{K}(\hat{\psi}_1^R, \hat{\psi}_{N_O}^O) \\ \vdots & \ddots & \vdots \\ \mathcal{K}(\hat{\psi}_{N_R}^R, \hat{\psi}_1^O) & \cdots & \mathcal{K}(\hat{\psi}_{N_R}^R, \hat{\psi}_{N_O}^O) \end{bmatrix}. \quad (6)$$

146 **2.3 Grasp Configuration Generation from  $\mathcal{D}(\mathcal{R}, \mathcal{O})$**

147 Given the predicted  $\mathcal{D}(\mathcal{R}, \mathcal{O})$ , we discuss how to generate the grasp joint values to grasp the object.  
 148 We first calculate the robot grasp point cloud, then estimate each link's 6D pose based on the joint  
 149 clouds. The system calculates the joint values by matching each link's 6D pose.

150 **Robotic Grasp Pose Point Cloud Generation** For a given point  $p_i^{\mathcal{R}}$ , the  $i$ -th row of  $\mathcal{D}(\mathcal{R}, \mathcal{O})$   
 151 denotes the distances from this robot grasp point to all points in the object point cloud. Given the  
 152 object point cloud, the multilateration method [21] positions the robot point cloud. This positioning  
 153 technique determines the location of a point  $p_i'^{\mathcal{R}}$  by solving the least-squares optimization problem  
 154 based on distances from multiple reference points:

$$p_i'^{\mathcal{R}} = \arg \min_{p_i'^{\mathcal{R}}} \sum_{j=1}^{N_{\mathcal{O}}} \left( \|p_i'^{\mathcal{R}} - p_j^{\mathcal{O}}\|_2^2 - \mathcal{D}(\mathcal{R}, \mathcal{O})_{ij}^2 \right)^2. \quad (7)$$

155 As shown in Zhou [26], this problem has a closed-form solution, and by using the implementa-  
 156 tion from Eisner et al. [25], we can directly compute  $p_i'^{\mathcal{R}}$ . Repeating this process for each row of  
 157  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  yields the complete predicted robot point cloud  $\mathbf{P}^{\mathcal{P}}$  in the grasp pose. In 3D space, we  
 158 can determine a point's position by measuring its relative distances to just three other points. Our  
 159  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation provides  $N_{\mathcal{O}} (= 512)$  relative distances, enhancing robustness to prediction  
 160 errors.

161 **6D Pose Estimation of Links** Directly solving inverse kinematics and getting the joint values from  
 162 a point cloud is not a trivial task. We first compute the 6D pose of each link in the world frame.  
 163 As described in Sec. 2.1, we store the point cloud for each link,  $\{\mathbf{P}_{\ell_i}\}_{i=1}^{N_{\ell}}$ . Given the predicted  
 164 grasp point cloud  $\{\mathbf{P}_{\ell_i}^{\mathcal{P}}\}_{i=1}^{N_{\ell}}$ , we calculate the 6D pose of each link using rigid body registration  
 165 techniques:

$$\mathcal{T}^* = (\mathbf{x}_i^*, \mathbf{R}_i^*) = \arg \min_{(\mathbf{x}_i, \mathbf{R}_i)} \|\mathbf{P}_{\ell_i}^{\mathcal{P}} - \mathbf{P}_{\ell_i}(\mathbf{x}_i, \mathbf{R}_i)\|^2, \quad (8)$$

166 where  $\mathbf{x}_i$  and  $\mathbf{R}_i$  represent the translation and rotation of the  $i$ -th link, respectively.

167 **Joint Configuration Optimization** After predicting the 6D pose for each link, our objective is to  
 168 optimize the joint values to align the translation of each link with the predicted result. Starting from  
 169 an initial value  $q_{init}$ , we iteratively solve the following optimization problem using CVXPY [27]:

$$\min_{\delta q} \left( \sum_{i=1}^{N_{\ell}} \left\| \mathbf{x}_i + \frac{\partial \mathbf{x}_i(q)}{\partial q} \delta q - \mathbf{x}_i^* \right\|_2 \right) \quad \text{s.t. } q + \delta q \in [q_{min}, q_{max}], |\delta q| \leq \varepsilon_q. \quad (9)$$

170 In each iteration, the system computes the delta joint values  $\delta q$  by minimizing the objective func-  
 171 tion and updates the joint values as  $q \leftarrow q + \delta q$ . Here,  $\mathbf{x}_i$  represents the current link translation,  
 172  $[q_{min}, q_{max}]$  denotes the joint limits, and  $\varepsilon_q = 0.5$  is the maximum allowable step size. The optimi-  
 173 zation process can be efficiently parallelized, typically achieving convergence within one second,  
 174 even for a 6+22 DoF ShadowHand.

175 **2.4 Loss Function**

176 The training objectives of the whole network include four parts, including the prediction of  $\mathcal{D}(\mathcal{R}, \mathcal{O})$   
 177 and  $\mathcal{T}$ , the suppression of penetration, and the KL divergence of the CVAE latent variable:

$$\begin{aligned} \mathcal{L} &= \lambda_{\mathcal{D}} \mathcal{L}_{\text{L1}} \left( \mathcal{D}(\mathcal{R}, \mathcal{O}), \mathcal{D}(\mathcal{R}, \mathcal{O})^{\text{GT}} \right) + \lambda_{\mathcal{T}} \frac{1}{N_{\ell}} \sum_{i=1}^{N_{\ell}} \mathcal{L}_{\ell_i} \\ &\quad + \lambda_{\mathcal{P}} |\mathcal{L}_{\mathcal{P}}(\mathbf{P}^{\mathcal{T}}, \mathbf{P}^{\mathcal{O}})| + \lambda_{KL} \mathcal{D}_{KL} \left( f_{\theta_G}(\mathbf{P}^G, \psi^{\mathcal{R}}, \psi^{\mathcal{O}}) \parallel \mathcal{N}(0, I) \right), \end{aligned} \quad (10)$$

178 where  $\lambda_{\mathcal{D}}$ ,  $\lambda_{\mathcal{T}}$ ,  $\lambda_{\mathcal{P}}$ ,  $\lambda_{KL}$  are hyperparameters for loss weights. The superscript 'GT' refers to the  
 179 ground truth annotations.  $\mathcal{N}(0, I)$  is a standard Gaussian distribution, and  $\mathbf{P}^{\mathcal{T}}$  is the robot point

180 cloud under the  $\mathcal{T}^*$  described in 2.3.  $\mathcal{L}_P$  computes the sum of the negative values of the signed  
 181 distance function (SDF) of  $\mathbf{P}^\mathcal{T}$  to  $\mathbf{P}^\mathcal{O}$  to penalize any penetration between the robot hand and the  
 182 object, and  $\mathcal{L}_\ell$  computes the difference between two 6D poses:

$$\mathcal{L}_{\ell_i} = \|\mathbf{x}_i^* - \mathbf{x}_i^{\text{GT}}\|_2 + \arccos\left(\frac{\text{tr}(\mathbf{R}_i^{*\text{T}} \mathbf{R}_i^{\text{GT}}) - 1}{2}\right). \quad (11)$$

183 Notably, the computation from  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  representation to the 6D pose  $\mathcal{T}^*$  shown in Eqn. 8 is en-  
 184 tirely matrix-based, ensuring differentiability for loss backpropagation and computational efficiency.

### 185 3 Experiments

186 In this section, we perform a series of experiments aimed at addressing the following questions  
 187 (Q1-Q6):

- 188 Q1: How successful are our generated grasps?
- 189 Q2: Does our unified model train on multi-embodiment outperform models trained on single em-  
 190 bodiments?
- 191 Q3: How diverse are our generated grasps?
- 192 Q4: How well does our pretraining learn configuration-invariant representations, and can this be  
 193 transferred across different embodiments?
- 194 Q5: How robust is our approach with partial object point cloud input?
- 195 Q6: How does our method perform in real-world settings?

#### 196 3.1 Evaluation Metric

197 **Success Rate:** We evaluate the success of grasping by determining whether the force closure con-  
 198 dition is satisfied. To implement this evaluation criterion, we used the Isaac Gym simulator [28]. A  
 199 simple PD controller is applied to execute the predicted grasps in the simulation. Certain forces are  
 200 applied sequentially along six orthogonal directions, following the approach in Li et al. [12]. We  
 201 apply each force for a duration of 1 second. We consider the grasp successful if the object’s resultant  
 202 displacement stays below 2 cm after applying the six directional forces.

203 **Diversity:** Grasp diversity is quantified by calculating the standard deviation of the joint values  
 204 (including 6 floating wrist DoF) across all successful grasps.

205 **Efficiency:** The computational time required to achieve a grasp is measured, encompassing both  
 206 network inference and the subsequent optimization steps.

#### 207 3.2 Dataset

208 We utilized a subset of the MultiDex dataset [12] (See Appendix for the filtering process). After  
 209 filtering, 24,764 valid grasps remained. We adopt three robots from the dataset: Barrett (3-finger),  
 210 Allegro (4-finger), and ShadowHand (5-finger). Each grasp defines its associated object, robot, and  
 211 grasp configurations. We retain the same training and test dataset splits as in the MultiDex dataset.

#### 212 3.3 Overall Performance

213 **Baselines** To answer Q1, we present a detailed comparison of  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  against DFC [2], GenDex-  
 214 Grasp [12], and ManiFM [13], as shown in Tab. 2. This comparison includes diverse methods to



Figure 4: Visualization of all methods.

Method	Success Rate (%) ↑				Diversity (rad.) ↑			Efficiency (sec.) ↓		
	Barrett	Allegro	ShadowHand	Avg.	Barrett	Allegro	ShadowHand	Barrett	Allegro	ShadowHand
DFC [2]	86.30	76.21	58.80	73.77	<b>0.532</b>	<b>0.454</b>	0.435	>1800	>1800	>1800
GenDexGrasp [12]	67.00	51.00	54.20	57.40	0.488	0.389	0.318	14.67	25.10	19.34
ManiFM [13]	-	42.60	-	42.60	-	0.288	-	-	9.07	-
DRO-Grasp (w/o pretrain)	87.20	82.70	46.70	72.20	<b>0.532</b>	0.448	0.429	<b>0.49</b>	<b>0.47</b>	<b>0.98</b>
<b>DRO-Grasp (Ours)</b>	<b>87.30</b>	<b>92.30</b>	<b>83.00</b>	<b>87.53</b>	0.513	0.397	<b>0.441</b>	<b>0.49</b>	<b>0.47</b>	<b>0.98</b>

Table 2: Overall comparison with baselines.

Method	Success Rate (%) ↑			Diversity (rad) ↑		
	Barrett	Allegro	ShadowHand	Barrett	Allegro	ShadowHand
Single	84.80	88.70	75.80	0.505	<b>0.435</b>	0.425
Multi	<b>87.30</b>	<b>92.30</b>	<b>83.00</b>	<b>0.513</b>	0.397	<b>0.441</b>
Partial	84.70	87.60	81.80	0.511	0.401	0.412

Table 3: Comparison under different conditions. “Single” trains on one hand, “Multi” trains on all hands, and “Partial” trains and tests on partial point clouds.

- 215 address the challenge of cross-embodiment grasping from various perspectives. They were evaluated  
 216 on 10 previously unseen test objects using the Barrett, Allegro, and ShadowHand robotic hands.  
 217 DFC is an optimization-based approach that searches for feasible grasp configurations through it-  
 218 erative optimization. GenDexGrasp predicts contact heatmaps and uses optimization to determine  
 219 grasp poses. ManiFM supports cross-embodiment grasping but employs a point-contact approach,  
 220 which was not suitable for training on our dataset that emphasizes surface-contact methods. As a  
 221 result, we can only evaluate its pretrained model of Allegro Hand for ManiFM.  
 222 Our experiments demonstrate that  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  significantly outperformed all baselines regarding suc-  
 223 cess rate across the robots by a large margin, highlighting the effectiveness of our approach. Fig. 4  
 224 visualizes grasps generated by our method alongside typical failure grasp poses from baselines. Our  
 225 approach generates reasonable grasps, while DFC often results in unnatural poses. GenDexGrasp  
 226 struggles with objects of complex shapes, frequently encountering significant penetration issues.  
 227 Although ManiFM produces visually appealing grasps, its point-contact method lacks stability, low-  
 228 ering its success rate in simulation.  
 229 From the first two rows of Tab. 3, we can see a slight improvement in success rates when training  
 230 across multiple robots compared to training on a single hand, demonstrating the cross-embodiment  
 231 generalizability of our method (Q2).  
 232 Our method significantly improves grasp generation speed. While DFC is slow in producing results  
 233 and learning-based methods like GenDexGrasp and ManiFM take tens of seconds per grasp due to  
 234 their complex optimization processes, our approach can generate a grasp within 1 second. This fast  
 235 computation is crucial for dexterous manipulation tasks.

### 236 3.4 Diverse Grasp Synthesis

237 Grasping diversity includes two key aspects: the wrist pose and the finger joint values. Since the  
 238 input and grasp rotations in the training data are correspondingly aligned, the model learns to implictly  
 239 map these rotations. This alignment enables the model, during inference, to generate appropriate  
 240 grasps based on the specified input orientation. Fig. 5 illustrates the grasp results for six different  
 241 input directions, showing that our model consistently produces feasible grasps, demonstrating the  
 242 controllability of our method. Additionally, by sampling the latent variable  $z \in \mathbb{R}^{64}$  from  $\mathcal{N}(0, I)$ ,

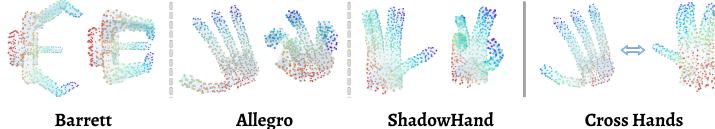


Figure 5: Diverse and pose-controllable grasp generation. The arrow refers to the input palm orientation. Arrows and hands of the same color represent corresponding input-output pairs.

243 our model can generate multiple grasps in the same direction, addressing Q3. As shown in Tab. 2,  
 244 the diversity of our method is highly competitive.

### 245 3.5 Configuration Correspondence Learning

246 As described in Sec. 2.1, our proposed configuration-invariant pretraining method learns an inherent  
 247 alignment across varying robotic hand configurations. To answer Q4, we visualize the learned cor-  
 248 respondence in Fig. 6, where each point in the closed-hand pose is colored according to the highest  
 249 cosine similarity with its counterpart in the open-hand pose. The excellent color matching within the  
 250 same hand demonstrates that the pretrained encoder successfully captures this alignment. Further-  
 251 more, strong matching across different hands highlights the transferability of features. As shown  
 252 in Tab. 2, removing the pretraining parameters and training the robot encoder directly results in  
 253 performance degradation across robotic hands, confirming the effectiveness of the pretrained model.



254 Figure 6: Visualization of the pretrained point matching.

### 254 3.6 Grasping with Partial Object Point Cloud Input

255 A common challenge in real-world experiments is the noise and incompleteness of point clouds from  
 256 depth cameras. Object-centric methods that rely on full object visibility often suffer performance  
 257 degradation under such conditions. In contrast, the relative distance feature of  $\mathcal{D}(\mathcal{R}, \mathcal{O})$  allows our  
 258 method to infer the robot point cloud even from partial observation. We validated this approach by  
 259 conducting experiments, removing 50% of the object point cloud in a contiguous region during both  
 260 training and validation. This setup simulates the incomplete data commonly encountered in practice.  
 261 As shown in the third row of Tab. 3, even with partial point clouds, our model can successfully  
 262 predict feasible grasps (Q5), indicating robustness when faced with incomplete input.

### 263 3.7 Real-Robot Experiments

264 We conducted real-world experiments with a uFactory xArm6 robot,  
 265 equipped with the LEAP Hand [29] and the overhead Realsense D435  
 266 camera, as illustrated in Fig. 7. As shown in Tab. 4, our method  
 267 achieved an average success rate of **89%** across 10 novel objects,  
 268 showcasing its effectiveness in dexterous grasping and its general-  
 269 izability to previously unseen objects (Q6). For experiment videos,  
 270 please visit our website <https://drograsp.github.io/>.



263 Figure 7: Real-world experiment setting.

Apple	Bag	Brush	Cookie Box	Cube	Cup	Dinosaur	Duck	Tea Box	Toilet Cleaner
9/10	10/10	9/10	10/10	9/10	7/10	9/10	8/10	8/10	10/10

263 Table 4: Real-world experiment results on unseen objects.

## 271 4 Conclusion

272 This work presents a new method for improving dexterous grasping by introducing the  $\mathcal{D}(\mathcal{R}, \mathcal{O})$   
 273 representation, which captures the essential interaction between robotic hands and objects. Unlike  
 274 existing methods that rely heavily on either object or robot-specific representations, our approach  
 275 bridges the gap by using a unified framework that generalizes well across different robots and ob-  
 276 ject geometries. Additionally, our pretraining approach enhances the model’s capacity to adapt to  
 277 different hand configurations, making it suitable for a wide range of robotic systems. Experi-  
 278 mental results confirm that our method delivers notable improvements in success rates, diversity, and  
 279 computational efficiency.

280 **References**

- 281 [1] M. A. Roa and R. Suárez. Grasp quality measures: review and performance. *Autonomous*  
282 *robots*, 38:65–88, 2015.
- 283 [2] T. Liu, Z. Liu, Z. Jiao, Y. Zhu, and S.-C. Zhu. Synthesizing diverse and physically stable grasps  
284 with arbitrary hand structures using differentiable force closure estimator. *IEEE Robotics and*  
285 *Automation Letters*, 7(1):470–477, 2021.
- 286 [3] S. Chen, J. Bohg, and C. K. Liu. Springgrasp: An optimization pipeline for robust and com-  
287 pliant dexterous pre-grasp synthesis. *arXiv preprint arXiv:2404.13532*, 2024.
- 288 [4] A. Patel and S. Song. GET-Zero: Graph embodiment transformer for zero-shot embodiment  
289 generalization, 2024. URL <https://arxiv.org/abs/2407.15002>.
- 290 [5] S. Haldar, J. Pari, A. Rai, and L. Pinto. Teach a robot to fish: Versatile imitation from one  
291 minute of demonstrations. *arXiv preprint arXiv:2303.01497*, 2023.
- 292 [6] Y. Xu, W. Wan, J. Zhang, H. Liu, Z. Shan, H. Shen, R. Wang, H. Geng, Y. Weng, J. Chen, et al.  
293 Unidexgrasp: Universal robotic dexterous grasping via learning diverse proposal generation  
294 and goal-conditioned policy. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
295 and Pattern Recognition, pages 4737–4746, 2023.
- 296 [7] W. Xu, W. Guo, X. Shi, X. Sheng, and X. Zhu. Fast force-closure grasp synthesis with learning-  
297 based sampling. *IEEE Robotics and Automation Letters*, 8(7):4275–4282, 2023.
- 298 [8] W. Wan, H. Geng, Y. Liu, Z. Shan, Y. Yang, L. Yi, and H. Wang. Unidexgrasp++: Improving  
299 dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-  
300 specialist learning. In *Proceedings of the IEEE/CVF International Conference on Computer*  
301 *Vision*, pages 3891–3902, 2023.
- 302 [9] L. Shao, F. Ferreira, M. Jorda, V. Nambiar, J. Luo, E. Solowjow, J. A. Ojea, O. Khatib, and  
303 J. Bohg. Unigrasp: Learning a unified model to grasp with multifingered robotic hands. *IEEE*  
304 *Robotics and Automation Letters*, 5(2):2286–2293, 2020.
- 305 [10] M. Attarian, M. A. Asif, J. Liu, R. Hari, A. Garg, I. Gilitschenski, and J. Tompson. Geometry  
306 matching for multi-embodiment grasping. In *Conference on Robot Learning*, pages 1242–  
307 1256. PMLR, 2023.
- 308 [11] S. Li, Z. Li, K. Han, X. Li, Y. Xiong, and Z. Xie. An end-to-end spatial grasp prediction model  
309 for humanoid multi-fingered hand using deep network. In *2021 6th International Conference*  
310 *on Control, Robotics and Cybernetics (CRC)*, pages 130–136. IEEE, 2021.
- 311 [12] P. Li, T. Liu, Y. Li, Y. Geng, Y. Zhu, Y. Yang, and S. Huang. Gendexgrasp: Generalizable dex-  
312 terous grasping. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*,  
313 pages 8068–8074. IEEE, 2023.
- 314 [13] Z. Xu, C. Gao, Z. Liu, G. Yang, C. Tie, H. Zheng, H. Zhou, W. Peng, D. Wang, T. Chen,  
315 Z. Yu, and L. Shao. Manifoundation model for general-purpose robotic manipulation of contact  
316 synthesis with arbitrary objects and robots, 2024.
- 317 [14] D. Morrison, P. Corke, and J. Leitner. Closing the loop for robotic grasping: A real-time,  
318 generative grasp synthesis approach. *arXiv preprint arXiv:1804.05172*, 2018.
- 319 [15] J. Varley, J. Weisz, J. Weiss, and P. Allen. Generating multi-fingered robotic grasps via deep  
320 learning. In *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*,  
321 pages 4415–4420. IEEE, 2015.
- 322 [16] A. Wu, M. Guo, and C. K. Liu. Learning diverse and physically feasible dexterous grasps with  
323 generative model and bilevel optimization. *arXiv preprint arXiv:2207.00195*, 2022.

- 324 [17] A. Goyal, V. Blukis, J. Xu, Y. Guo, Y.-W. Chao, and D. Fox. Rvt2: Learning precise manipulation  
325 from few demonstrations. *RSS*, 2024.
- 326 [18] A. Goyal, J. Xu, Y. Guo, V. Blukis, Y.-W. Chao, and D. Fox. Rvt: Robotic view transformer  
327 for 3d object manipulation. *CoRL*, 2023.
- 328 [19] Y. Huang, C. Agia, J. Wu, T. Hermans, and J. Bohg. Points2plans: From point clouds to  
329 long-horizon plans with composable relational dynamics. *arXiv preprint arXiv:2408.14769*,  
330 2024.
- 331 [20] W. Huang, C. Wang, Y. Li, R. Zhang, and L. Fei-Fei. Rekep: Spatio-temporal reasoning of  
332 relational keypoint constraints for robotic manipulation. *arXiv preprint arXiv:2409.01652*,  
333 2024.
- 334 [21] A. Norrdine. An algebraic solution to the multilateration problem. In *Proceedings of the*  
335 *15th international conference on indoor positioning and indoor navigation, Sydney, Australia*,  
336 volume 1315, 2012.
- 337 [22] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon. Dynamic graph  
338 cnn for learning on point clouds. *ACM Transactions on Graphics (tog)*, 38(5):1–12, 2019.
- 339 [23] A. Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*,  
340 2017.
- 341 [24] K. Sohn, H. Lee, and X. Yan. Learning structured output representation using deep conditional  
342 generative models. *Advances in neural information processing systems*, 28, 2015.
- 343 [25] B. Eisner, Y. Yang, T. Davchev, M. Vecerik, J. Scholz, and D. Held. Deep se (3)-equivariant  
344 geometric reasoning for precise placement tasks. *arXiv preprint arXiv:2404.13478*, 2024.
- 345 [26] Y. Zhou. An efficient least-squares trilateration algorithm for mobile robot localization. In  
346 *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3474–  
347 3479. IEEE, 2009.
- 348 [27] S. Diamond and S. Boyd. CVXPY: A Python-embedded modeling language for convex opti-  
349 mization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- 350 [28] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox. Gpu-accelerated  
351 robotic simulation for distributed reinforcement learning. In *Conference on Robot Learning*,  
352 pages 270–282. PMLR, 2018.
- 353 [29] K. Shaw, A. Agarwal, and D. Pathak. Leap hand: Low-cost, efficient, and anthropomorphic  
354 hand for robot learning. *Robotics: Science and Systems (RSS)*, 2023.