

# Geometry Matching for Multi-Embodiment Grasping

Anonymous Author(s)

Affiliation

Address

email

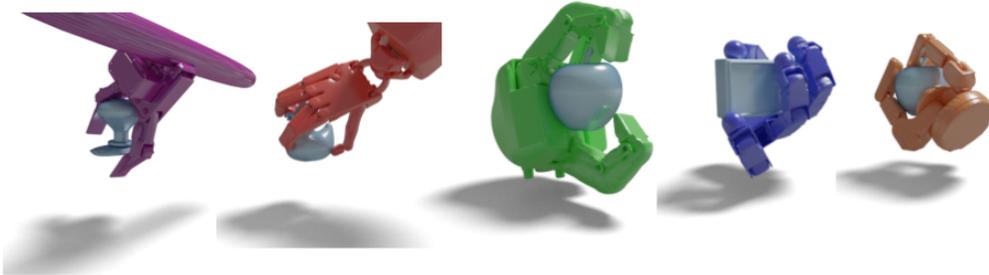


Figure 1: **GeoMatch**: Our method enables multi-embodiment grasping by conditioning the grasp selection on end-effector and object geometry.

1       **Abstract:** While significant progress has been made on the problem of generating  
2 grasps, many existing learning-based approaches still concentrate on a single em-  
3 bodiment, provide limited generalization to higher DoF end-effectors and cannot  
4 capture a diverse set of grasp modes. In this paper, we tackle the problem of grasping  
5 multi-embodiments through the viewpoint of learning rich geometric representa-  
6 tions for both objects and end-effectors using Graph Neural Networks (GNN).  
7 Our novel method - *GeoMatch* - applies supervised learning on grasping data from  
8 multiple embodiments, learning end-to-end contact point likelihood maps as well  
9 as conditional autoregressive prediction of grasps keypoint-by-keypoint. We com-  
10 pare our method against 3 baselines that provide multi-embodiment support. Our  
11 approach performs better across 3 end-effectors, while also providing competitive  
12 diversity of grasps. Examples can be found at [geo-match.github.io](https://github.com/geo-match).

13       **Keywords:** Multi-Embodiment, Dexterous Grasping, Graph Neural Networks

## 14   1 Introduction

15 Dexterous grasping remains an open and important problem for robotics manipulation. Many tasks  
16 where robots are involved, from the simplest to the most complex ones, at their core come down to  
17 some form of interacting with objects in their environment. This in turn, results in grasping objects  
18 with all kinds of different geometries. In addition, the large variety of robot and end-effector types  
19 necessitates that grasping should also be achievable with new and arbitrary end-effector geometries.  
20 However, the cross-embodiment gap between grippers does not permit simply applying grasping  
21 policies from one end-effector to another, while domain adaptation i.e. “translating” actions from  
22 one embodiment to another, is also not straightforward. In comparison, humans are extremely versa-  
23 tile: they can adapt the way they grasp objects based on what they know about object geometry even  
24 if the object class or instance is new to them, and they can do this in more than one ways efficiently.

25 There has been much research in grasping thus far, with many works focusing on one embodiment  
26 at a time [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11] and fewer looking at the multi-embodiment problem [12,  
27 13, 14]. Methods are divided between hand-agnostic or hand-aware, and experiment with different  
28 representations for grasping, such as contact maps [12], contact points [13] or even root pose and

29 z-offset [14]. Existing multi-embodiment approaches either require explicit representation of joint  
30 limits that becomes exponentially harder in higher DoF end-effectors, or expect heavy manual work  
31 to adapt to new end-effectors, or showcase mixed rates of success across different embodiments with  
32 some gripper/hand morphologies performing significantly better than others.

33 Drawing inspiration from how humans seem to be adapt their grasps easily and successfully based on  
34 priors they have learned about 3D geometry of both objects in space and their own hands, we propose  
35 endowing robotic agents with a similar sense of geometry via a generalized geometry embedding  
36 that is able to represent both objects and end-effectors. This embedding can be used to predict grasps  
37 that demonstrate stability, diversity, generalizability, and robustness. More specifically, we propose  
38 to rely on Graph Neural Networks (GNN) to encode meaningful geometry representations and use  
39 them to predict keypoint contacts on the object surface in an autoregressive manner.

40 In summary, our contributions are as follows:

- 41 a) We propose formulating robot grasp learning as a geometry matching problem through learning  
42 correspondences between geometry features, where both end-effector and object geometries are  
43 encoded in rich embedding spaces.
- 44 b) To solve the above problem formulation, we introduce a novel method, namely GeoMatch, that  
45 is trained end-to-end to learn expressive geometric embeddings, and autoregressive keypoint con-  
46 tacts via teacher forcing.
- 47 c) We demonstrate that our method is competitive against baselines without any extra requirements  
48 to support higher DoF end-effectors and while also showcasing high performance across multiple  
49 embodiments.

## 50 2 Related Work

51 **Dexterous Grasping.** Many dexterous grasping works do not look into multi-embodiment and in-  
52 stead focus on diversity of objects using a single end-effector. Many grasping methods support  
53 2-finger parallel grippers [2, 3, 4, 5, 6, 7, 8, 15] with several others looking into high-DoF dexterous  
54 manipulation [9, 10, 1, 16]. Some work has also been conducted towards multi-embodiment grasp-  
55 ing. Several of those address the problem from the differentiable simulation grasp synthesis point  
56 of view [17, 11, 18]. GenDexGrasp [12] advocate for hand-agnostic contact maps generated by a  
57 trained cVAE [19] with a newly introduced align distance, and optimize matching of end-effectors  
58 against produced contact maps via a specialized Adam optimization method. This is the most re-  
59 cent work to our knowledge, attempting to tackle multi-embodiment grasping without extra steps  
60 to support higher DoF grippers. In contrast to them, we choose to operate on hand-specific contact  
61 maps as we are interested in learning both object and embodiment geometry conditioned grasps,  
62 and empirically found our method to perform more evenly well multi-embodiments. Intuitively, our  
63 work is closest to UniGrasp [13]. UniGrasp operates on object and end-effector point clouds to  
64 extract features and ultimately output contact points which are then fed into an Inverse Kinematics  
65 (IK) solver, similarly to us. Their encoder of choice is PointNet++ and the contact prediction is done  
66 through a Point Set Selection Network (PSSN) [20]. Their proposed architecture adds one stage per  
67 finger, which means supporting more than 3 finger grippers requires manually adding another stage.  
68 As a result, adapting the method to more than 2-finger and 3-finger grippers requires significant  
69 work, while also the need for explicit representation of boundary configurations can explode expo-  
70 nentially on higher DoF end-effectors. In contrast, we rely on learned geometry features to identify  
71 viable configurations as opposed to explicitly encoding them through joint limit representation, as  
72 well as on a small number of user-selected keypoints, same for all end-effectors, which disentangles  
73 the dependency between number of fingers and applicability of our method. Similarly to UniGrasp,  
74 EfficientGrasp [21] also uses PointNet++ and a PSSN model for contact point prediction and further  
75 generates a pose with RL. TAX-Pose [22] is another recent work that shares some high level con-  
76 cepts. Instead of encoding the end-effector, authors look at the problem of tasks involving objects  
77 that interact with each other in a particular way. They proceed with encoder objects or object parts  
78 using DGCNN and learn a cross-attention model that predicts relative poses of objects that accom-  
79 plish a task. AdaGrasp [14] uses 3D Convolutional Neural Networks to learn a scoring function for  
80 possible generated grasps, and finally executes the best one. Many of the methods mentioned, rely  
81 on deterministic solvers which can result in decreased diversity of generated grasps. Even though  
82 we also rely on a deterministic solver, we address this issue by leveraging the scoring we obtain by

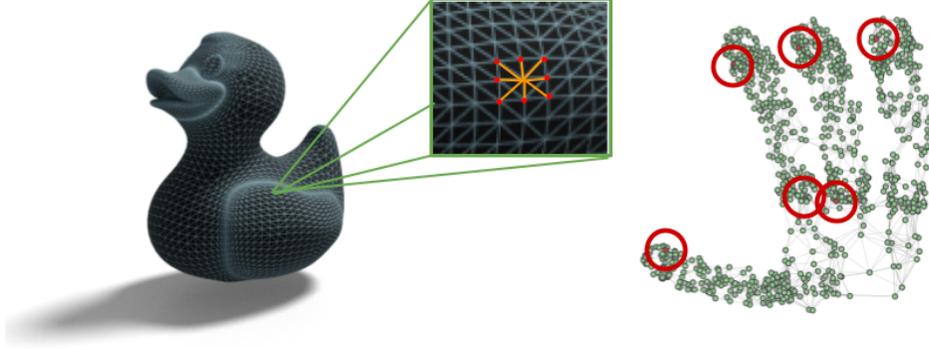


Figure 2: **Object and end-effector inputs.** Objects are initially represented as regularly sampled point clouds which are converted into a graph for further processing. End-effector geometries are given as meshes and converted to coarser graphs by randomly sampling points from the mesh as an intermediate step. User-selected keypoints are highlighted in red.

83 the learned full unnormalized distribution of contacts to select a first keypoint that will guide the  
 84 remaining contact point prediction. This permits higher diversity without having to sample a large  
 85 number of grasps.

86 **Graph Neural Networks.** Graph Neural Networks were first introduced by Scarselli et al. [23] as  
 87 a proposed framework to operate on structured graph data. Since then, many advancements have  
 88 been made towards extending their capabilities and expressivity [24]. Specifically in the grasping  
 89 literature, there have also been multiple instances of use of GNN. More specifically, Huang et al. [25]  
 90 propose learning a GNN to predict 3D stress and deformation fields based on finite element method  
 91 based grasp simulations. The use of GNNs for end-effector parameterization has been proposed  
 92 before in [26] where tactile sensor data is fed into a GNN to represent the end-effector as part of grasp  
 93 stability prediction, however we propose applying GNN as a more general geometry representation  
 94 that encompasses both objects and end-effectors jointly. Lou et al. [27] leverage GNN to represent  
 95 the spacial relation between objects in a scene and suggest optimal 6-DoF grasping poses. Unlike  
 96 previous methods, we aim to use GNN as a general geometry representation for any rigid body,  
 97 including both objects and end-effectors. For the purposes of this work, we leverage the GNN  
 98 implementation by [28] due to the readily available and easily adaptable code base.

99 **Geometry-Aware Grasping.** In the topic of geometry-aware grasping, several works have advo-  
 100 cated for the importance of geometry in the grasping problem. Yan et al. [29] encodes RGBD input  
 101 via generative 3D shape modeling and 3D reconstruction, then based on this learned geometry-  
 102 aware representation grasping outcomes are predicted with solutions coming out of an analysis-by-  
 103 synthesis optimization. In the same vein, Van et al. [30] proposed leveraging learned 3D reconstruc-  
 104 tion as a means of understanding geometry, and further rely on this for grasp success classification  
 105 as an auxiliary objective function for grasp optimization and boundary condition checking. Bohg et  
 106 al. [31] introduced a supervised learning method where a classifier trained on labeled images pre-  
 107 dicted grasps via shape context based representations. Finally, Jiang et al. [6] learn grasp affordances  
 108 and 3D reconstruction as an auxiliary task, through the use of implicit functions. Unlike these works,  
 109 we suggest looking at geometry itself directly from 3D as a feature representation without imposing  
 110 any 3D reconstruction constraints.

### 111 3 Method

112 In this work, we aim to learn robust and performant grasping prediction via embeddings of geometry  
 113 for both objects and end-effectors. We are given point cloud representations of object and end-  
 114 effector geometries. These are converted into graphs which allows to utilize GNNs to learn features  
 115 across both.

116 Assume an object geometry represented as a graph  $\mathcal{G}_O = (\mathcal{V}_O, \mathcal{E}_O)$  and an end-effector geometry  
 117 also represented as a graph  $\mathcal{G}_G = (\mathcal{V}_G, \mathcal{E}_G)$  where  $\mathcal{V}_O, \mathcal{V}_G, \mathcal{E}_O, \mathcal{E}_G$  are the object and end-effector  
 118 vertices and edges respectively. The edges are represented by adjacency matrices  $Adj_O, Adj_G$  for the  
 119 object and end-effector graphs respectively, which are row normalized symmetric binary matrices

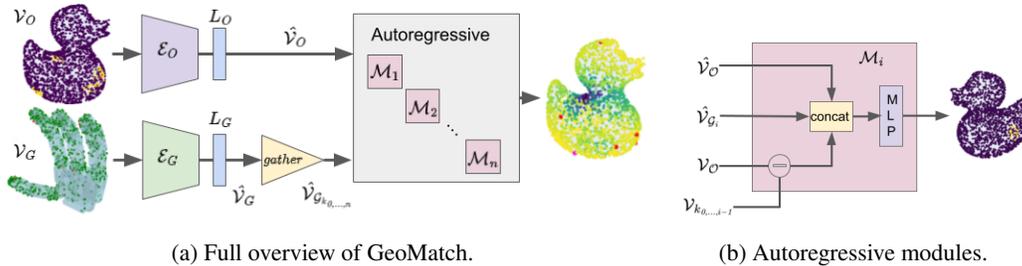


Figure 3: **GeoMatch architecture.** The object and gripper graphs are passed through the two encoders followed by linear layers. The gripper keypoint embeddings are gathered and are passed as input along with the object embeddings in the autoregressive modules.

120 with a unitary diagonal. Given  $\mathcal{G}_O$  and  $\mathcal{G}_G$ , we seek to learn feasible and stable contact points  
 121 between a subset of  $\mathcal{V}_O$  and  $\mathcal{V}_G$ .

### 122 3.1 Object and End-effector Representations

123 Each end-effector is represented by its surface geometry in the form of a graph. Additionally, we  
 124 require a small number of canonical user-selected keypoints that will be the ones matched with object  
 125 vertices when calculating contacts. We select these once for each end-effector we are working with  
 126 visually, and store them. It is recommended to selected keypoints having good coverage of each  
 127 gripper with respect to its morphology and its grasping behavior. To construct the graph, we sample  
 128 a point cloud of 1000 surface points from the end-effector mesh in a canonical rest pose:

$$q_{\text{rest}} = (t, R, \theta_0, \dots, \theta_{N-1})_{\text{rest}}, \quad (1)$$

129 where  $t \in \mathbb{R}^3$  is the root translation,  $R \in \mathbb{R}^6$  is the root rotation in continuous 6D representation  
 130 as introduced in [32], and  $\theta_0, \dots, \theta_{N-1}$  are the joint angles of the end-effector. We chose the rest  
 131 pose to be a vector with all joint angles in middle range of their respective joint limits, zero root  
 132 translation, and identity root rotation. For creating the graph, we consider each of the points a vertex  
 133 and create edges between each point and its  $K$  closest points.

134 For our experiments, we empirically chose  $K = 8$  to capture local geometry. This is a hyperparameter  
 135 that depends on point cloud density and object structure. The canonical keypoints were selected  
 136 manually at the rest pose of each hand to represent points of contact. We empirically chose 6 so that  
 137 for all end-effectors in our dataset, each finger and the palm is represented by at least one keypoint.  
 138 Choosing the same number of keypoints for different embodiments is technically not a requirement  
 139 as lower degree of freedom end-effectors may have good enough coverage with less, but we chose to  
 140 use a constant number multi-embodiments to simplify our training process. A sample of the object  
 141 and end-effector representations is shown in Fig. 2.

142 Each object is also represented by its surface geometry in the form of a graph. More specifically, the  
 143 same process is utilized to convert a object point cloud of 2048 points to a graph. The point cloud  
 144 and adjacency matrix together describe the graph. For the purposes of this work, we use a subset of  
 145 the MultiDex dataset introduced by [12] and used by them to train the CMap-VAE model of their  
 146 approach. The dataset is comprised of 5 end-effectors - one 2-finger, two 3-finger, one 4-finger, and  
 147 one 5-finger, as well as 58 common objects from YCB [33] and ContactDB [34]. It contains 50,802  
 148 diverse grasps over the set of hands and objects, each represented by an object name, an end-effector  
 149 name, and the end-effector pre-grasp pose in the form of Eq. 1.

### 150 3.2 Learning Setup

151 At the core of our hypothesis is that learning rich geometry features for objects and end-effectors  
 152 jointly can be a powerful tool for dexterous and diverse grasp prediction multi-embodiments. Thus,  
 153 we seek an architecture that can embed local geometry information well. From the architecture  
 154 choices that demonstrate such properties, we chose Graph Neural Neural Networks (GNN) [28].

155 Our overall model architecture can be seen in Fig. 3a and is designed to learn a) an independent un-  
 156 normalized prior distribution of contacts between each object and gripper keypoint, and b) marginal

157 distributions of contact for each keypoint, conditioned upon the above likelihood map and previously  
 158 predicted keypoint contacts.

### 159 3.2.1 Independent unnormalized priors of object-keypoint contacts

160 Ideally, what we would like to calculate is the full joint distribution of vertex-to-vertex contact, i.e.  
 161  $P(v_i, v_j)$  for all vertices of the object  $v_i$ , and all vertices of the end-effector  $v_j$ . This is typically  
 162 intractable. We reduce the complexity by only focusing on  $n$  landmark keypoints  $k_0, \dots, k_n$  on the  
 163 end-effector and we try to approximate  $P(v_o, k_0, \dots, k_n)$  through learning a set of factorizations by  
 164 applying the Bayes rule. This yields

$$P(v_o, k_0, \dots, k_n) = \prod_{i=1}^n P_{M_i}(v_o, k_i | k_0, \dots, k_{i-1}) = \prod_{i=1}^n P_{M_i}(v_o, k_i | \mathbf{k}_{<i}), \quad (2)$$

165 where  $v_o \in \mathcal{V}_O$ ,  $(k_0, \dots, k_n) \subset \mathcal{V}_G$ , and  $P_{M_i}(v_o, k_i | k_0, \dots, k_{i-1})$  are the factorized marginals to  
 166 be learned in an autoregressive manner, as discussed in the following subsection. As a first step,  
 167 we aim to associate a likelihood of contact for a sparse set of keypoints  $k_i$  per each object vertex  
 168  $v_o$ . We first pass the object and end-effector graphs through GNN encoders that output the same  
 169 number of features. The embeddings obtained are L2 normalized. We then gather the embeddings  
 170 on the canonical user-selected keypoints as the vertices of interest on the hand. It is noted that  
 171 we still compute the embedding for all hand vertices even though for contact areas, we focus on  
 172 the embedding of the canonical keypoints. This unnormalized likelihood map of object-keypoint  
 173 contacts intuitively represents a score that a given object vertex is in contact with a given gripper  
 174 keypoint and is given by

$$P_{I_i}(v_o, k_i) = E_O(v_o) \cdot E_G(v_g)[k_i]. \quad (3)$$

175 This is optimized against the dot product of the hand-specific object contact map  $C_O(v_o, k_i)$  via a  
 176 binary cross-entropy loss

$$\mathcal{L}_{P_{I_0, \dots, n}} = \sum_{i=1}^n \text{BCE}_{\lambda_a}(P_{I_i}(v_o, k_i), C_O(v_o, k_i)), \quad (4)$$

177 where  $\lambda_a$  is the positive weight hyperparameter used to address the class imbalance.

### 178 3.2.2 Autoregressive marginals with teacher forcing

179 As discussed in the previous paragraph, we seek to estimate the joint distribution of contacts by  
 180 estimating a set of factorizations. We further proceed with the estimation of factors:

$$P_{M_i}(v_o, k_i | k_0, \dots, k_{i-1}) \quad \forall i \in [0, n]. \quad (5)$$

181 Both, object and gripper embeddings are projected down to a lower dimension with a simple linear  
 182 layer without bias, and passed into 5 layers, each responsible for predicting the index of the object  
 183 vertex  $v_{o_n}$  where keypoint  $k_n$  makes contact, given keypoints  $k_{0 \dots (n-1)}$ .

184 Each layer  $n$  concatenates the embedding of the  $n$ -th keypoint of the end-effector along with the  
 185 object embedding. Then, it calculates the *relative* distance map of each object vertex to each of the  
 186  $n - 1$  object vertices where the previous  $n - 1$  keypoints make contact. Note that is done via teacher  
 187 forcing: instead of using the predictions of each  $n - 1$  layer, we use the previous  $n - 1$  ground  
 188 truth contact points. This avoids error propagation during training. The relative distance maps are  
 189 stacked and concatenated with the object and  $n$ -th keypoint embeddings. This constitutes the input  
 190 to an MLP that predicts a binary classification prediction over the object vertices that indicates the  
 191 predicted  $n$ -th contact point.

192 This is again optimized against the ground truth binary contact map label of the  $n$ -th gripper key-  
 193 point, contributing to a second binary cross-entropy loss term

$$\mathcal{L}_{P_{M_0, \dots, n}} = \sum_{i=1}^n \text{BCE}_{\lambda_b}(P_{M_i}(v_o, k_i | k_0, \dots, k_{i-1}), C_O(v_o, k_i)), \quad (6)$$

194 where, similarly,  $\lambda_b$  is the positive weight hyperparameter used to address the class imbalance. A  
 195 visual representation of the autoregressive layers can be seen in 3b. Note that for  $i = 0$ ,  $P_{I_0}(v_o, k_0)$   
 196 constitutes the first marginal for  $k_0$  and thus:  $P_{I_0}(v_o, k_0) = P_{M_0}(v_o, k_0)$ .

## 197 3.3 Likelihood Maps

198 In order to learn the above, we assumed access to ground truth likelihood maps used for supervised  
 199 learning which we obtain as follows. For each grasp in our dataset, instead of an object contact map,

200 we generate a (2048, 6) per-gripper-keypoint proximity map where the nearest areas are calculated  
 201 as a fixed number of  $M$  closest points in Euclidean distance, to each of the canonical keypoints:

$$P_o(v_o, k_i) = \begin{cases} 1, & v_o \in \arg \min_M \|\mathcal{V}_O - \mathcal{V}_G(k_i)\|^2 \text{ for each } k_i, \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

202 We also generate a gripper contact map for the selected keypoints where the contacts are defined as  
 203 the keypoints closer than a given threshold, to the object point cloud:

$$C_g(k_i) = \begin{cases} 1, & \exists v_o, \|\mathcal{V}_O - \mathcal{V}_G(k_i)\|^2 < \text{threshold}, \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

204 where  $P_o(v_o, k_i)$  is the object proximity map,  $C_g(k_i)$  is the gripper contact map,  $O(v_o)$  is the object  
 205 point on index  $v_o$  and  $G(k_i)$  is the gripper point on the canonical keypoint index  $k_i$ . For this work,  
 206 we empirically assumed  $M = 20$  and a threshold of 0.04. Finally, the hand-specific object contact  
 207 map can be obtained as  $C_O(v_o, k_i) = P_o(v_o, k_i) \cdot C_g(k_i)$ .

208 For training speed considerations, we preprocess the dataset prior to training, and save each grasp in  
 209 this new form.

210 We finally represent our full training objective with the total loss being:

$$\mathcal{L}_{\text{total}} = \alpha \cdot \mathcal{L}_{P_{I_0, \dots, n}} + \beta \cdot \mathcal{L}_{P_{M_0, \dots, n}} \quad (9)$$

211 For our experiments, we used the Graph Convolutional Networks (GCN) implementation by Kipf et  
 212 al. [28] with 3 hidden layers of size 256,  $\lambda_a = 500$ ,  $\lambda_b = 200$ , 512 output embedding dimension,  
 213 one for objects and one for end-effectors. The linear projection for each encoder was of size 64  
 214 without bias. We also used  $\alpha = \beta = 0.5$ .

### 215 3.4 Grasp Prediction at Inference

216 At test time, the independent unnormalized distribution for  $k = 0$  is leveraged to sample keypoint  
 217 0 which will commence the autoregressive inference. More specifically, we use the 0-th dimension  
 218 as a scoring mechanism for sampling high likelihood points where keypoint 0 makes contact. This  
 219 is then passed into the model as the previous contact of keypoint 1. At inference time, at the  $n$ -th  
 220 step, teacher forcing is substituted with passing in the  $(n - 1)$  predicted contact vertices. Finally,  
 221 the end result is a tensor of 6 coordinates of the object graph. As previously mentioned, grasping  
 222 is a multi-modal distribution and our model should be able to sample from the various modes. In  
 223 our method, this can be achieved straightforwardly by sampling a variety of starting top-K points  
 224 for keypoint 0. The intuition behind this is that diverse, yet likely starting points for keypoint 0 will  
 225 condition subsequent predicted points differently, and ultimately yield different grasp modes. For  
 226 our experiments, we sampled 4 such top-K points, namely top-0, 20, 50, 100, in order to explore the  
 227 capacity of our method to generate diverse grasps. A more sophisticated sampling algorithm such as  
 228 Beam search, could be applied here, however we empirically achieved sufficient diversity through  
 229 multimodal sampling of keypoint 0. Here it should be noted that this autoregressive representation  
 230 does present some limitations. More specifically, the ordering with which the keypoint contacts are  
 231 being learning and ultimately selected could vastly change the result. However, we refrain from  
 232 experimenting with all possible combinations of keypoint ordering in the scope of this work.

233 The end-effector joint angles are then inferred by feeding the predicted contact points into an In-  
 234 verse Kinematics (IK) solver. For our purposes, we used SciPy’s Trust Region Reflective algo-  
 235 rithm (TRF) [35]. The initial pose given to IK is a heuristic pose calculated by applying a rota-  
 236 tion/translation that aligns the palm with the closest object vertex while keeping all non-root joints  
 237 at their rest pose configuration. It should be noted that any other IK solution and initial pose guess  
 238 strategy could be leveraged instead. Further implementation details can be found in Appendix B.

## 239 4 Experiments

240 We evaluate our method through the lens of a number of research questions.

241 **Q1: How successful is the model at producing stable and diverse grasps for various embodi-  
 242 ments?** We train our method with a training set containing samples of 5 end-effectors and 38 objects.



Figure 4: **Qualitative results.** Generated grasps using GeoMatch on **unseen** objects with ezgripper, barrett, robotiq-3finger, allegro and shadowhand. For each grasp, another perspective is included where the GeoMatch predicted keypoints on each object are marked with purple and the gripper user selected keypoints matching these, are marked with yellow.

243 We then generate grasps on each of the 5 end-effectors but 10 new unseen objects, and evaluate them  
 244 in Isaac Gym [36], specifically the Isaac Gym based environment proposed by [12]. Similarly, we  
 245 apply a consistent  $0.5ms^{-2}$  acceleration on the object from all  $xyz$  directions sequentially, for 1  
 246 second each. If the object moves more than  $2cm$  after every such application, the grasp is declared  
 247 as a failure. We also follow the same contact-aware refinement paradigm, which applies force clo-  
 248 sure via a single step of Adam with step size 0.05. In addition, we provide calculated diversity as  
 249 the standard deviation of the joint angles of all successful grasps, comparably to [12]. We compare  
 250 our method to GenDexGrasp [12], AdaGrasp (initOnly as it is the closest setup to our task) [14], and  
 251 DFC [17]. Results can be found in Tab. 1. In addition, we provide a number of qualitative results in  
 252 Fig. 4.

Method	Success (%) $\uparrow$				Diversity (rad) $\uparrow$		
	ezgripper	barrett	shadowhand	Mean	ezgripper	barrett	shadowhand
DFC [17]	58.81	85.48	72.86	72.38	0.3095	0.3770	0.3472
AdaGrasp [14]	60.0	80.0	-	70.0	0.0003	0.0002	-
GenDexGrasp [12]	43.44	71.72	<b>77.03</b>	64.01	0.238	0.248	0.211
<b>GeoMatch (Ours)</b>	<b>75.0</b>	<b>90.0</b>	72.5	<b>79.17</b>	0.188	0.249	0.205

Table 1: **Success and diversity comparisons.** GeoMatch performs more evenly well across end-effectors with a varied DoF number while maintaining diversity of grasp configurations.

253 In our experiments, we observed that GeoMatch is performing slightly worse (-2%) on the 5-finger  
 254 gripper Shadowhand than the best performing baseline, however performance for the 2-finger and  
 255 3-finger grippers **increases by 5-30%** compared to other methods. Diversity remains competitive to  
 256 other methods. Overall, the minimum performance observed for GeoMatch is significantly higher  
 257 than baselines and the average performance multi-embodiments beats all baselines we compared  
 258 against.

259 **Q2: Is the multi-embodiment model performing better than a model trained on individual em-  
 260 bodiments?** We hypothesize that training our method on data containing a variety of end-effectors  
 261 will result in learning better geometry representations. To investigate this, we train our method on  
 262 each single embodiment separately by filtering our dataset for each given end-effector. We then com-  
 263 pare against the multi-embodiment model. Each of the single end-effector models is trained only on  
 264 grasp instances of that gripper while the multi-embodiment model is trained on all 5 end-effectors  
 265 and objects in the training set. The validation set in all cases contains 10 unseen objects. We provide  
 266 results in Tab. 2. The model trained on multi-embodiment data is indeed performing **20%-35%**  
 267 better than single end-effector models which advocates for the value of multi-embodiment grasping  
 268 policies as opposed to single model policies trained on more data.

269 **Q3: How robust is the learned model under relaxed assumptions?** While our method demon-  
 270 strates compelling results, it has been trained on full point clouds. Acknowledging that this is often  
 271 a strict assumption, especially when considering real-world environments, we evaluate robustness  
 272 of the approach under conditions more similar to real-world robotic data. We experiment with grasp  
 273 generation using: a) noisy point clouds, b) partial point clouds, and c) partial point clouds including

Method	Success (%) $\uparrow$			Diversity (rad) $\uparrow$		
	ezgripper	barrett	shadowhand	ezgripper	barrett	shadowhand
Single embodiment	40.0	70.0	40.0	0.157	0.175	0.154
Multi embodiment	<b>75.0</b>	<b>90.0</b>	<b>72.5</b>	0.188	0.249	0.205

Table 2: Comparisons between the Multi-embodiment model and models trained on individual grippers.

274 noise. For each of these, we perturbed the object point clouds accordingly, and collected grasps  
 275 using our method **zero-shot**. Success rate was **77.5%**, **66.7%**, and **67.5%** across end-effectors for  
 276 each type of augmentation respectively. As demonstrated, our method shows reasonable robustness.  
 277 Experiment details and a breakdown of numbers can be found at Appendix A.

278 **Q4: How important are various components of the design?** Finally, we investigate the design  
 279 decisions of our approach and how they affect performance. More specifically, we perform two  
 280 ablations:

281 PointNet++ as the encoder of choice instead of GNN. We evaluate our choice towards GNN by  
 282 swapping out the two GNN encoders with PointNet++[20], a popular encoder architecture for point  
 283 clouds. Our results show that GNN was indeed a good choice as it performs better than the Point-  
 284 Net++ ablation, by **10%** averaging across end-effectors. In addition, we empirically observed a 12x  
 285 slow down when using PointNet++ due to the difference in model parameters number, which also  
 286 makes GNN more light weight and fast. A breakdown per end-effector can be found in Appendix A.

287 Non-shared weights between keypoint encoders. We hypothesize that a shared encoder among all  
 288 end-effectors is beneficial for learning features that represent local geometry and this subsequently,  
 289 informs autoregressive prediction of keypoints. To validate this hypothesis, we conducted an abla-  
 290 tion where we separated the end-effector encoder to 6 separate identical encoders, one per keypoint.  
 291 Our main model with shared weights across all end-effectors and keypoints outperforms the split  
 292 encoders by **9%**. Further analysis per end-effector can be found in Appendix A.

## 293 5 Limitations

294 While this method showcased that grasp learning can benefit from multi-embodiment data in terms  
 295 of generalization to new objects as well as robustness, obtaining large amounts of such multi-  
 296 embodiment grasping data, especially in real world setups can be challenging, time consuming and  
 297 expensive. However, given that a single embodiment grasping policy was shown to require more  
 298 data to perform comparably, we argue that spending resources on a multi-embodiment dataset to  
 299 yield a policy that performs well across a variety of grippers is a better choice. Lastly, our method  
 300 relies on the robustness of the IK solution. We empirically observed cases where there was a rea-  
 301 sonable grasp solution for a set of predicted keypoints, however the chosen IK solution terminated  
 302 in some suboptimal configuration.

## 303 6 Conclusion

304 This work presented a novel multi-embodiment grasping method that leverages GNN to learn pow-  
 305 erful geometry features for object and embodiment representation. Our approach demonstrates that a  
 306 joint encoder trained on multiple embodiments can better embed geometry in a generalizable fashion  
 307 and ultimately result in higher grasping success rate on unseen objects. The proposed framework  
 308 also showcased robustness to more realistic point cloud inputs. Diversity of generated grasps re-  
 309 mains competitive while producing such diverse grasps is as simple as conditioning with a different  
 310 high likelihood starting contact point for the first keypoint. Code and models will be released on  
 311 acceptance.

## References

- 312
- 313 [1] A. Nguyen, D. Kanoulas, D. G. Caldwell, and N. G. Tsagarakis. Detecting Object Affordances  
314 with Convolutional Neural Networks. In *International Conference on Intelligent Robots and*  
315 *Systems (IROS)*, 2016.
- 316 [2] H.-S. Fang, C. Wang, M. Gou, and C. Lu. GraspNet-1Billion: A Large-Scale Benchmark for  
317 General Object Grasping. In *Conference on Computer Vision and Pattern Recognition (CVPR)*,  
318 2020.
- 319 [3] A. Mousavian, C. Eppner, and D. Fox. 6-dof graspnet: Variational grasp generation for object  
320 manipulation. In *International Conference on Computer Vision (ICCV)*, 2019.
- 321 [4] A. ten Pas, M. Gualtieri, K. Saenko, and R. Platt. Grasp pose detection in point clouds. *The*  
322 *International Journal of Robotics Research*, 36(13-14), 2017.
- 323 [5] M. Gualtieri, A. ten Pas, K. Saenko, and R. Platt. High precision grasp pose detection in dense  
324 clutter. In *International Conference on Intelligent Robots and Systems (IROS)*, 2016.
- 325 [6] Z. Jiang, Y. Zhu, M. Svetlik, K. Fang, and Y. Zhu. Synergies between affordance and geometry:  
326 6-dof grasp detection via implicit representations. *arXiv preprint arXiv:2104.01542*, 2021.
- 327 [7] S. Wang, Z. Zhou, and Z. Kan. When transformer meets robotic grasping: Exploits context for  
328 efficient grasp detection. *Robotics and Automation Letters*, 7(3), 2022.
- 329 [8] C. Wu, J. Chen, Q. Cao, J. Zhang, Y. Tai, L. Sun, and K. Jia. Grasp Proposal Networks:  
330 An End-to-End Solution for Visual Learning of Robotic Grasps. In *Conference on Neural*  
331 *Information Processing Systems (NeurIPS)*, 2020.
- 332 [9] Y. Xu, W. Wan, J. Zhang, H. Liu, Z. Shan, H. Shen, R. Wang, H. Geng, Y. Weng, J. Chen, et al.  
333 Unidexgrasp: Universal robotic dexterous grasping via learning diverse proposal generation  
334 and goal-conditioned policy. In *Proceedings of the Conference on Computer Vision and Pattern*  
335 *Recognition*, 2023.
- 336 [10] W. Wan, H. Geng, Y. Liu, Z. Shan, Y. Yang, L. Yi, and H. Wang. UniDexGrasp++: Improving  
337 Dexterous Grasping Policy Learning via Geometry-aware Curriculum and Iterative Generalist-  
338 Specialist Learning. *arXiv preprint arXiv:2304.00464*, 2023.
- 339 [11] D. Turpin, L. Wang, E. Heiden, Y.-C. Chen, M. Macklin, S. Tsogkas, S. Dickinson, and  
340 A. Garg. Grasp’d: Differentiable contact-rich grasp synthesis for multi-fingered hands. In  
341 *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27,*  
342 *2022, Proceedings, Part VI*, 2022.
- 343 [12] P. Li, T. Liu, Y. Li, Y. Geng, Y. Zhu, Y. Yang, and S. Huang. Gendexgrasp: Generalizable  
344 dexterous grasping. *arXiv preprint arXiv:2210.00722*, 2022.
- 345 [13] L. Shao, F. Ferreira, M. Jorda, V. Nambiar, J. Luo, E. Solowjow, J. A. Ojea, O. Khatib,  
346 and J. Bohg. Unigrasp: Learning a unified model to grasp with multifingered robotic hands.  
347 *Robotics and Automation Letters*, 5(2), 2020.
- 348 [14] Z. Xu, B. Qi, S. Agrawal, and S. Song. Adagrasp: Learning an adaptive gripper-aware grasping  
349 policy. In *International Conference on Robotics and Automation (ICRA)*, 2021.
- 350 [15] M. Sundermeyer, A. Mousavian, R. Triebel, and D. Fox. Contact-graspnet: Efficient 6-dof  
351 grasp generation in cluttered scenes. In *International Conference on Robotics and Automation*  
352 *(ICRA)*, 2021.
- 353 [16] S. Brahmabhatt, A. Handa, J. Hays, and D. Fox. ContactGrasp: Functional Multi-finger Grasp  
354 Synthesis from Contact. In *International Conference on Intelligent Robots and Systems (IROS)*,  
355 2019.
- 356 [17] T. Liu, Z. Liu, Z. Jiao, Y. Zhu, and S.-C. Zhu. Synthesizing diverse and physically stable  
357 grasps with arbitrary hand structures using differentiable force closure estimator. *Robotics and*  
358 *Automation Letters*, 7(1), 2021.

- 359 [18] D. Turpin, T. Zhong, S. Zhang, G. Zhu, E. Heiden, M. Macklin, S. Tsogkas, S. Dickinson,  
360 and A. Garg. Fast-Grasp'D: Dexterous Multi-finger Grasp Generation Through Differentiable  
361 Simulation. In *ICRA*, 2023.
- 362 [19] K. Sohn, H. Lee, and X. Yan. Learning Structured Output Representation using Deep Condi-  
363 tional Generative Models. In *Advances in Neural Information Processing Systems*, 2015.
- 364 [20] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on  
365 point sets in a metric space. *Advances in neural information processing systems*, 2017.
- 366 [21] K. Li, N. Baron, X. Zhang, and N. Rojas. Efficientgrasp: A unified data-efficient learning to  
367 grasp method for multi-fingered robot hands. *Robotics and Automation Letters*, 7(4), 2022.
- 368 [22] C. Pan, B. Okorn, H. Zhang, B. Eisner, and D. Held. TAX-Pose: Task-Specific Cross-Pose  
369 Estimation for Robot Manipulation. In *Conference on Robot Learning (CoRL)*, 2022.
- 370 [23] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The Graph Neural  
371 Network Model. *Transactions on Neural Networks*, 2009.
- 372 [24] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun. Graph neural  
373 networks: A review of methods and applications. *AI open*, 1, 2020.
- 374 [25] I. Huang, Y. Narang, R. Bajcsy, F. Ramos, T. Hermans, and D. Fox. DefGraspNets: Grasp  
375 Planning on 3D Fields with Graph Neural Nets. *arXiv preprint arXiv:2303.16138*, 2023.
- 376 [26] A. Garcia-Garcia, B. S. Zapata-Impata, S. Orts-Escolano, P. Gil, and J. Garcia-Rodriguez.  
377 Tactilegn: A graph convolutional network for predicting grasp stability with tactile sensors.  
378 In *International Joint Conference on Neural Networks (IJCNN)*, 2019.
- 379 [27] X. Lou, Y. Yang, and C. Choi. Learning object relations with graph neural networks for target-  
380 driven grasping in dense clutter. In *International Conference on Robotics and Automation*  
381 *(ICRA)*, 2022.
- 382 [28] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks.  
383 *arXiv preprint arXiv:1609.02907*, 2016.
- 384 [29] X. Yan, J. Hsu, M. Khansari, Y. Bai, A. Pathak, A. Gupta, J. Davidson, and H. Lee. Learning  
385 6-DOF Grasping Interaction via Deep Geometry-Aware 3D Representations. In *International*  
386 *Conference on Robotics and Automation (ICRA)*, 2018.
- 387 [30] M. Van der Merwe, Q. Lu, B. Sundaralingam, M. Matak, and T. Hermans. Learning continuous  
388 3d reconstructions for geometrically aware grasping. In *International Conference on Robotics*  
389 *and Automation (ICRA)*, 2020.
- 390 [31] J. Bohg and D. Kragic. Learning grasping points with shape context. *Robotics and Autonomous*  
391 *Systems*, 58(4), 2010.
- 392 [32] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li. On the continuity of rotation representations in  
393 neural networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- 394 [33] B. Calli, A. Singh, J. Bruce, A. Walsman, K. Konolige, S. Srinivasa, P. Abbeel, and A. M. Dol-  
395 lar. Yale-CMU-Berkeley dataset for robotic manipulation research. *The International Journal*  
396 *of Robotics Research*, 36(3), 2017.
- 397 [34] S. Brahmhatt, C. Ham, C. C. Kemp, and J. Hays. ContactDB: Analyzing and Predicting Grasp  
398 Contact via Thermal Imaging. In *Conference on Computer Vision and Pattern Recognition*  
399 *(CVPR)*, 2019.
- 400 [35] J. J. Moré and D. C. Sorensen. Computing a Trust Region Step. *SIAM Journal on Scientific*  
401 *and Statistical Computing*, 4(3), 1983.
- 402 [36] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin,  
403 A. Allshire, A. Handa, and G. State. Isaac Gym: High Performance GPU-Based Physics  
404 Simulation For Robot Learning, 2021.