Neural Grasp Distance Fields for Robot Manipulation

Thomas Weng^{1,2}, David Held², Franziska Meier¹, and Mustafa Mukadam¹

¹Meta AI, ²Carnegie Mellon University

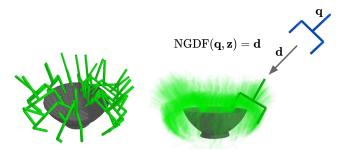
Abstract—We formulate grasp learning as a neural field and present Neural Grasp Distance Fields (NGDF). Here, the input is a 6D pose of a robot end effector and output is a distance to a continuous manifold of valid grasps for an object. In contrast to current approaches that predict a set of discrete candidate grasps, the distance-based NGDF representation is easily interpreted as a cost, and minimizing this cost produces a successful grasp pose. This grasp distance cost can be incorporated directly into a trajectory optimizer for joint optimization with other costs such as trajectory smoothness and collision avoidance. During optimization, as the various costs are balanced and minimized, the grasp target is allowed to smoothly vary, as the learned grasp field is continuous. We evaluate NGDF on joint grasp and motion planning in simulation and the real world, outperforming baselines by 63% execution success while generalizing to unseen query poses and unseen object shapes. Project page: https://sites. google.com/view/neural-grasp-distance-fields.

I. INTRODUCTION

We present *Neural Grasp Distance Fields* (NGDF), which model the continuous manifold of valid grasp poses as the level set of a neural implicit function. Given a 6D query pose, NGDF predicts the unsigned distance between the query and the closest valid grasp on the manifold (see Fig. 1).

Neural implicit fields have driven recent advancements in novel view synthesis [1] and 3D reconstruction [2], [3], [4], [5]. These approaches represent distributions as continuous functions that take a query as input and predict its relationship to the learned distribution. In 3D shape reconstruction, for instance, neural implicit fields are used to represent the surface of a shape: 3D points are used as queries, and the output is the distance to the surface, or occupancy at the query point. Unlike explicit methods, neural implicit fields can encode complex topological distributions and are not limited by resolution.

With NGDF, formulating grasp learning as a neural field allows us to interpret the implicit function as a cost such that a query pose can be optimized to result in a grasp pose. Prior grasp estimation methods largely output a discrete set of candidate grasps [6], [7], [8], [9], [10], from which one grasp must be selected to perform downstream planning. Instead, we incorporate the grasp distance cost directly into a gradient-based optimizer [11] to jointly optimize the grasp and reaching motion from an initial trajectory. During each optimization iteration, NGDF estimates the distance between the final gripper pose of the trajectory and the grasp level set. This "grasp distance" is minimized as a cost, along with other trajectory costs such as smoothness and collision avoidance. The gradient of the grasp cost for updating the trajectory is computed through fully differentiable operations. This



(a) Discrete Grasp Set

(b) Continuous Grasp Manifold with NGDF

Fig. 1: (a) Existing grasp estimation methods produce discrete grasp sets which do not represent the true continuous manifold of possible grasps. (b) Our work, Neural Grasp Distance Fields (NGDF), learns a continuous grasp manifold: given a query pose **q** and an object shape embedding **z**, NGDF outputs the distance **d** between **q** and the closest grasp. This distance can be leveraged as a cost for optimization, facilitating joint grasp and motion planning.

optimization results in a smooth, collision-free trajectory that reaches a valid grasp pose.

In experiments, we find that NGDF learns the level set of valid grasp poses, outperforms baselines by 63% execution success on simulated reaching and grasping, and generalizes to unseen object shapes and poses in the real world. The key contributions of the paper are:

- Neural Grasp Distance Fields (NGDF), a neural implicit function that predicts the distance between a query pose and the closest grasp, representing the manifold of grasps as a continuous level set.
- A gradient-based optimization algorithm that incorporates NGDF for joint reach and grasp planning.

II. RELATED WORK

While grasping and motion planning are well-studied topics in robotics, prior works often propose different system designs with different assumptions, making comparison and contextualization difficult. We summarize the most important design decisions for 6-DOF grasp and motion planning and trace the decisions in representative methods (see Fig. 2).

A. 6-DOF Grasp Estimation

6-DOF grasp estimation is a well-studied task [17], [18] that aims to predict successful grasps in SE(3) for target objects; we focus here on recent, data-driven methods. State-of-the-art methods take point clouds as input and output a discrete set of grasps, representing only a subset of the true continuous grasp set [7], [6], [8], [9], [10]. Outputting a finer discretization comes with a cost of a greater computational

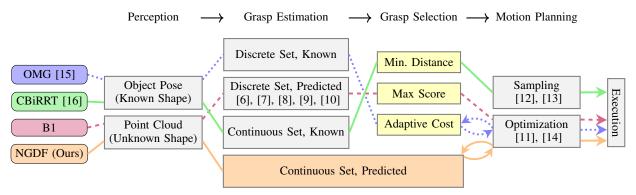


Fig. 2: Columns illustrate design decisions within grasp and motion planning pipelines. The left-most column highlights representative pipelines like OMG-Planner [15], CBiRRT [16], and baseline B1 from Table II which uses a SOTA grasp estimator [6]. The respective design choices for these methods are traced through the columns. We identify learned continuous representations as an under-explored option for grasp estimation, and propose NGDF as a solution that does not require a heuristic grasp selection step since the grasp pose is jointly optimized with motion planning.

complexity for both grasp estimation as well as grasp selection: a final grasp must be chosen from the predicted set. Because these methods only predict discrete grasp sets, they necessitate a multi-stage approach, which can be brittle if any of the stages (grasp estimation, selection, or motion planning) fails. Our single-stage approach models grasps as the level set of a continuous implicit function to jointly optimize grasping and motion planning.

B. Joint Grasp Selection and Motion Planning

Following the multi-stage paradigm above, several works assume a grasp set is provided by an upstream method, and address the downstream task of planning a reaching trajectory. Berenson et al. [16] model grasp sets as a continuous range of poses called Task Space Regions, and use samplingbased planning to satisfy the constraint. GOMP [19], uses sequential quadratic programming on discrete grasp sets for fast bin picking. Goal-set CHOMP [20] incorporates hard constraints like goal sets into trajectory optimization. The methods above do not address the problem of switching between grasps during planning; OMG-Planner [15] therefore proposes online learning to estimate goal costs and switch to the minimum cost grasp at every optimization iteration. OMG-Planner used ground-truth grasp sets per object, though their method can use estimated grasp sets as well. Our approach does not assume grasps are provided and does not require explicit grasp selection; instead, NGDF estimates and updates the grasp pose during trajectory optimization itself.

Other works propose closed-loop methods for 6-DOF grasping. Wang *et al.* [21] learn a latent space of trajectories for closed-loop grasping. Song *et al.* [22] learn a closed-loop policy from human demonstrations. Temporal Grasp-Net [23] updates a discrete grasp set over time by querying a grasp evaluator. In this work, we introduce a novel implicit representation for the grasp manifold. We focus on open-loop planning and leave closed-loop planning with NGDF as future work.

C. Implicit Neural Representations

Recent advances in vision and graphics research have used implicit neural representations to achieve impressive results

on novel view synthesis [1] and 3D reconstruction [3], [5], [2], [4]. Karunratakul *et al.* [24] learn an implicit representation for human grasp poses. Inspired by these works, we learn an implicit neural function to predict distances between query gripper poses and grasp poses, and use this function to optimize grasp trajectories.

The robotics community has also explored neural implicit functions for a variety of manipulation tasks [25], [26], [27], [28], [29], [30], [31]. GIGA [32] proposed using neural implicit functions to model both 3D shape and grasp quality. However, GIGA predicts a single grasp parameterization per 3D location, and requires a sampling procedure to select the final pose from the implicit set. Our approach predicts grasp distance, allowing multiple grasp orientations per 3D location, and uses optimization to minimize grasp distance and achieve the grasp pose.

Concurrent works have proposed continuous representations for dexterous hands [33] and multiple grippers [34]. Urain *et al.* [35] represents grasps as diffusion fields, framing joint grasp and motion planning as an inverse diffusion process. In this paper, we use an implicit function to represent grasp distance, and use gradient-based trajectory optimization for joint grasp and motion planning.

III. BACKGROUND

Neural Implicit Functions. Neural implicit functions (NIFs) are neural networks that take a query $\mathbf{q} \in \mathbb{R}^d$ and optionally a context embedding $\mathbf{z} \in \mathcal{Z}$ to output a scalar value that represents a relationship to an underlying distribution: $f(\mathbf{q}, \mathbf{z}) : \mathbb{R}^d \times \mathcal{Z} \mapsto \mathbb{R}$. In the domain of 3D shape reconstruction, the context \mathbf{z} is a latent shape embedding, the query \mathbf{q} is a 3D point, and the scalar output is either distance to the closest surface [2], [4], or occupancy [3], [5]. The shape surface is represented by the zero level set in distance-based methods, or the decision boundary in occupancy-based methods. Unlike explicit functions, NIFs are not limited by resolution as they predict a value at any query point, and also better represent underlying distributions that are disjoint [28]. Our approach leverages both properties in learning a manifold of grasps.

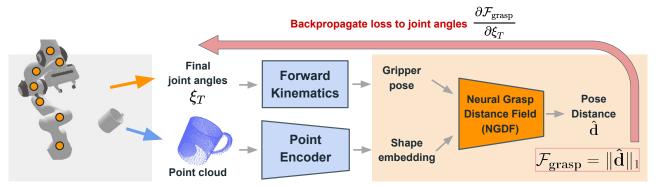


Fig. 3: We use NGDF as a goal cost function on the final state of a trajectory during gradient-based optimization. Given the current robot joint configuration and a point cloud of an object or scene, the current gripper pose and a shape embedding are computed as inputs for NGDF. Then, NGDF predicts the distance of the current gripper pose to the closest grasp (Sec. IV-A). The predicted distance is used as the cost and the gradient with respect to the joint configuration is computed with backpropagation. This cost (with gradient) is used with other costs like smoothness and collision avoidance to update the trajectory (Sec. IV-B).

Gradient-based Trajectory Optimization. A mapping from time t to robot joint configuration \mathbf{p} is defined as a trajectory $\xi:[0,T]\to\mathbf{p}$. Trajectory optimization aims to find the optimal trajectory given an objective functional \mathcal{U} :

$$\xi^* = \arg\min_{\xi} \mathcal{U}[\xi], \text{ s.t. } \xi(0) = \mathbf{p}_s, \ \xi(T) = \mathbf{p}_g$$
 (1)

for a given start \mathbf{p}_s and goal \mathbf{p}_g configuration. In manipulation, the objective \mathcal{U} contains cost terms for smoothness and collision avoidance. CHOMP [11] solves for ξ^* with functional gradient descent:

$$\xi_{t+1} = \xi_t - \eta A^{-1} \bar{\nabla} \mathcal{U}(\xi_t) \tag{2}$$

where A is an acceleration metric that helps propagate updates over the entire trajectory.

IV. METHOD

In this work, we represent a set of poses $\mathcal{M} \subset \mathbf{SE}(3)$ as the level set of a neural implicit function. This implicit function takes a query pose \mathbf{q} as input and estimates its distance to the learned level set. Sec. IV-A describes how Neural Grasp Distance Fields (NGDF) leverage this insight to learn the level set of valid *grasp* poses. Sec. IV-B explains how to incorporate NGDF into a trajectory optimization framework to jointly reason over smooth and collision-free reaching trajectories that end at a valid grasp pose. Fig. 3 provides an overview of our method.

A. Neural Grasp Distance Fields

Given a query pose $\mathbf{q} \in \mathbf{SE}(3)$ and a shape embedding $\mathbf{z} \in \mathcal{Z}$, NGDF defines an implicit function: NGDF $(\mathbf{q}, \mathbf{z}) = \mathbf{d}$, where \mathbf{d} is the distance from \mathbf{q} to the closest valid grasp $\mathbf{g} \in \mathcal{M} \subset \mathbf{SE}(3)$ for an object in a scene. Valid grasps are poses where a gripper can stably grasp an object by closing its fingers. For the distance metric \mathbf{d} we combine translation and orientation distances into a single "control points" metric [7]:

$$d_i = \|\mathcal{T}(\mathbf{q}; \mathbf{c}_i) - \mathcal{T}(\mathbf{g}; \mathbf{c}_i)\|_1, \ i = 0, \dots, N$$
 (3)

where $\mathcal{T}(\cdot; \mathbf{c}_i)$ is the transformation of a predefined set of points $\{\mathbf{c}_i\}$ on the gripper. Since \mathbf{q} and \mathbf{g} belong to $\mathbf{SE}(3)$, the distance could be defined based on the manifold geodesic distance between those poses, however we find that the

control points based distance metric balances the translation and rotation costs better in practice. NGDF estimates the distance for each control point $\mathbf{c}_{0...N}$ separately: $\mathbf{d}(\mathbf{q}, \mathbf{g}) = \begin{bmatrix} d_0, \ldots, d_N \end{bmatrix}^T$. During training, the estimated distances $\hat{\mathbf{d}}$ are supervised with L1 loss: $\mathcal{L} = \|\hat{\mathbf{d}} - \mathbf{d}\|_1$.

B. Optimization of Grasping Trajectories using NGDF

For a given query pose, NGDF outputs the distance to the closest grasp pose. We now show how to enable joint optimization for reaching and grasping with NGDF. We incorporate NGDF as a goal cost estimator within a gradientbased trajectory optimizer that already has cost terms for smoothness and collision avoidance.

In this work, we combine NGDF with CHOMP [11] (described in Sec. III), though NGDF can be used in any gradient-based trajectory optimization algorithm. Since CHOMP specifies a fixed goal \mathbf{p}_g , we modify CHOMP to include \mathbf{p}_g as a variable in the optimization following Dragan *et al.* [20]. We then add our grasp cost $\mathcal{F}_{\text{grasp}}$ as the variable goal cost to the objective functional \mathcal{U} :

$$\mathcal{U}[\xi] = \lambda_1 \mathcal{F}_{\mathbf{grasp}}[\xi] + \lambda_2 \mathcal{F}_{\mathrm{smooth}}[\xi] + \lambda_3 \mathcal{F}_{\mathrm{obs}}[\xi]$$
 (4)

where λ_i are cost weights.

Grasp Distance as a Goal Cost. We now define $\mathcal{F}_{\mathrm{grasp}}$ and derive its functional gradient $\bar{\nabla}\mathcal{F}_{\mathrm{grasp}}$ for gradient-based optimization. For a trajectory (during any iteration of optimization), we calculate the gripper pose from the final joint configuration using forward kinematics: $\mathbf{q}_T = \mathrm{FK}(\xi_T)$. We then use NGDF to estimate the distance of this gripper pose to a valid grasp: $\mathrm{NGDF}(\mathbf{q}_T,\mathbf{z}) = \hat{\mathbf{d}}$. The norm of this distance becomes our grasp cost: $\mathcal{F}_{\mathrm{grasp}}[\xi] = \|\hat{\mathbf{d}}\|_1$. We can compute the gradient of the grasp cost with respect to the joint configuration ξ_T through backpropagation:

$$\frac{\partial \mathcal{F}_{\text{grasp}}}{\partial \xi_T} = \frac{\partial \mathcal{F}_{\text{grasp}}}{\partial \mathbf{q}_T} \frac{\partial \mathbf{q}_T}{\partial \text{FK}} \frac{\partial \text{FK}}{\partial \xi_T}$$
 (5)

Since the grasp cost only applies to the final configuration in a trajectory, the functional gradient $\bar{\nabla} \mathcal{F}_{\mathrm{grasp}}$ contains all zeros except for the last row: $\bar{\nabla} \mathcal{F}_{\mathrm{grasp}} = [\mathbf{0}, \mathbf{0}, \dots, \frac{\partial \mathcal{F}_{\mathrm{grasp}}}{\partial \xi_T}]^T$.

Joint Optimization of Trajectory Costs. Similar to the objective functional (Eq. 4), the objective functional gradient $\bar{\nabla}\mathcal{U}$ is a weighted sum of gradients: $\bar{\nabla}\mathcal{U}[\xi] = \lambda_1 \bar{\nabla}\mathcal{F}_{\text{grasp}} + \lambda_2 \bar{\nabla}\mathcal{F}_{\text{smooth}} + \lambda_3 \bar{\nabla}\mathcal{F}_{\text{obs}}$. At every optimization iteration, we compute the costs and functional gradients as described above, then update the trajectory according to the A-metric update rule (Eq. 2). Since our objective cost has terms for minimizing distance to a valid grasp, maintaining smoothness, and avoiding collisions, our algorithm jointly optimizes all three to produce reaching and grasping trajectories.

C. Implementation Details

Dataset. Training NGDF requires a dataset of point clouds, valid grasp poses, and query poses. We use the ACRONYM [36] dataset, which contains object meshes and successful grasp poses collected in NVIDIA FleX [37]. For grasp poses, our evaluations in Sec. V are run in PyBullet [38], so we relabel the successful grasp poses based on their success in PyBullet with the same linear and rotational shaking parameters used in ACRONYM. In addition, we filter the positive grasp set to only include grasps where the normals at the mesh and finger contact points are opposed to each other (-0.98 cosine similarity). Our results in Sec. V-B show that this filtering improves grasp performance. To collect query poses for the dataset, we sample 1 million random SE(3) poses within a 0.5 m radius of the object mesh centroid. While it is possible that some of the sampled poses could be positive grasps, we assume they are few in number and do not run additional grasp evaluation to filter them. For each sampled pose, we use distance to the closest grasp in the valid grasp set (see Sec. IV-A) as our supervision.

Architecture. An input point cloud is converted into the shape embedding **z** using a VN-OccNet [39] encoder pretrained on 3D reconstruction [29]. The input to NGDF is a concatenation of this shape embedding **z** with the input query **q**'s position and quaternion. The NGDF network is based on DeepSDF [2] and consists of 8 MLP layers, 512 units each, and ReLU activations on the hidden layers. A softplus activation on the output layer ensures positive outputs.

Training Procedure. We freeze the weights of the pretrained point encoder during training and only train the NGDF network. Each training sample consists of a partial point cloud, a query pose, and the closest valid grasp. Similar to NDF [29], the partial point cloud is merged together from 4 camera views and downsampled to 1500 points using farthest point sampling. Random rotation augmentations are applied to each sample with 70% probability. Finding the ground truth closest grasp pose is computationally expensive and requires multiple simulated grasp attempts per query pose. Therefore, our supervision is pseudo-ground truth, as the closest grasp pose comes from a large but discrete set of grasps [36]. We find that this discrete grasp set is dense enough to train NGDF, while still representing unseen valid grasp poses at or near the zero level set (Sec. V-A).

Trajectory Optimization. CHOMP [11] uses a fixed or decaying step size for functional gradient updates, which is sufficient for trajectories with fixed start and goal joint

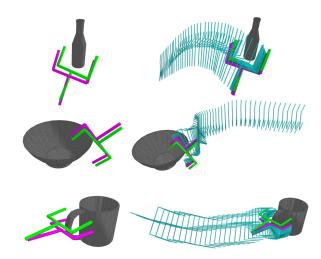


Fig. 4: Grasp Level Set Evaluation. Left: Final predicted pose (magenta) and its closest grasp pose (green) in the training dataset. Right: Gripper path (teal) as it is optimized from initial to final pose. Object meshes shown for visual clarity; our method takes point clouds as input.

configurations. However, with our modification of CHOMP in Sec. IV-B to allow a variable goal configuration, we found that such simple step size strategies resulted in poor convergence. We address this issue by using Adam [40] to adaptively update the step size ("CHOMP-Adam"). We use differentiable SE(3) operations [41] and a differentiable robot model [42] to backpropagate gradients from the output of NGDF to the robot joint configuration (Eq. 5).

V. EXPERIMENTS

We first evaluate how well NGDFs represent valid grasp manifolds as their zero level sets (Sec. V-A). Then we perform a full system evaluation with NGDFs on a "reaching and grasping" task (Sec. V-B), where an NGDF is used within a gradient-based trajectory optimizer as a goal cost function. We evaluate generalization on grasping intracategory unseen objects (Sec. V-C), and demonstrate grasping on a real robot system (Sec. V-D).

A. NGDF Level Set Evaluation

First, we investigate whether the learned level set of an NGDF represents successful grasps. Our evaluation procedure considers driving an initial query pose to the learned level set. We use the distance output from NGDF as a loss, and update the query pose with Adam [40] using backpropagated gradients. Note that this evaluation optimizes just the gripper pose; full-arm trajectory optimization is considered in the next subsection. We evaluate NGDF on three objects: Bottle, Bowl, and Mug. For this evaluation, we train a single NGDF model for each object, and evaluate models trained with and without the dataset filtering procedure described in Sec. IV-C. We run the optimization for 3k steps with a learning rate of 1e-4. Since we represent poses as positions and quaternions, we normalize the quaternion after each gradient update to ensure valid rotations.

TABLE I: NGDF Grasp Level Set Results

	Train Set Error (m) ↓	Grasp Success ↑
Bottle-NoFilter Bottle	0.023 ± 0.01 0.029 ± 0.01	0.480 0.880
Bowl-NoFilter	0.036 ± 0.02	0.540
Bowl	0.033 ± 0.01	0.760
Mug-NoFilter	0.038 ± 0.01	0.680
Mug	0.035 ± 0.01	0.860

Results are averaged over 50 unseen query poses per object, sampled from within a 0.5 m radius of the object centroid.

The quantitative results on grasp level set optimization are shown in Table I. We use two metrics for this evaluation. The "Train Set Error" metric is the minimum control points distance (Eq. 3) between the optimized gripper pose and the closest grasp pose in the discrete training set. Since NGDF should learn a continuous level set and interpolate between grasps in the training set, we expect NGDF not to achieve zero error on this metric, but it provides a good surrogate for comparing models. The "Grasp Success" metric measures the grasp quality of the optimized gripper poses. For each pose, we load the target object in PyBullet [38] and attempt a grasp at the specified pose. The robot gripper is always initialized to the same position; the object is transformed relative to the gripper. Linear and rotational shaking are applied after gripping the object [36], and the grasp is successful if the object is still gripped after the shaking.

Our results show that while NGDFs trained on filtered and unfiltered data have similar Train Set Error, the Grasp Success for filtered data models is much higher. These results also indicate that NGDFs have learned continuous level sets, since the mean distance predicted by NGDF after optimization is less than 1e-5, much lower than the minimum distance to the training set of grasps. Fig. 4 shows examples of the optimization path and achieved gripper pose.

B. Simulated Reaching and Grasping Evaluation

Next, we evaluate our method on a full reaching and grasping task, which requires planning a smooth, collision-free grasping trajectory for the full robot arm starting from an initial robot joint configuration. This evaluates the full pipeline as opposed to just the stand-alone gripper pose in the previous subsection. The task is considered successful if the robot executes the trajectory, closes its fingers to grasp the object, and lifts the object without losing it. We place Bottle, Bowl, and Mug objects in simulation in 30 random orientations each (see Appendix Fig. S2 in [43], left-most column), thus 90 trials in total. Our results indicate that even in a seemingly simple setting, randomly oriented objects present an overall challenging benchmark.

For this evaluation, we train a separate NGDF (similar to NeRF approaches [1], [26]) for each object, though our method can be extended to generalize across objects like other shape-conditioned implicit approaches [2]. We also evaluate intra-category (known class, unseen shape) generalization in the next subsection. We run 500 iterations of CHOMP-Adam (see Sec. IV-C) with a learning rate of

3e-3. The grasp cost is weighted heavily relative to the collision and smoothness costs. The trajectory is initialized using inverse kinematics so the gripper pose of the final joint configuration is within 0.3 m of the center of the object point cloud; the rest of the initial trajectory is interpolated between the start and end joint configurations.

The results are shown in Table II. We compare against oracle methods that provide upper-bound task performance, and against baselines that predict discrete grasps. Oracle methods assume perfect object pose estimation and known discrete grasp set. All discrete grasp methods run inverse kinematics over all discrete grasp goals and discard infeasible grasps. For planning, methods use goal-set CHOMP [20] or CHOMP [11], depending on whether the goal is fixed or can vary. "O1" selects the goal with minimum distance to the initial joint configuration, and keeps it fixed throughout planning. "OMG" [15] adaptively learns a cost for each grasp and selects the grasp with minimum cost at every optimization iteration (Variable Goal).

The baselines that predict discrete grasps use Contact-Graspnet [6] as the grasp estimator. We use weights (provided by the authors) that are trained on millions of grasps and shapes. "B1" selects the grasp goal with the maximum score estimated by Contact-GraspNet and keeps it fixed during planning. "B2" selects the grasp goal with minimum distance to the initial joints and keeps it fixed during planning. "B3" allows varying grasps during planning using the minimum distance metric. "B4" uses the same adaptive cost from OMG [6] to select grasp goals during planning.

Our results show that while oracle methods perform well, methods that don't assume known object pose and use predicted grasps have much lower Execution Success. Of the predicted grasp methods, NGDF performs best. Surprisingly, the B3 and B4 variable goal variants do not outperform fixed goal variants B1 and B2. Failure cases for all methods are largely due to collisions between the gripper fingers and the object, which are a relatively small obstacle cost and may be difficult for the planner to balance with the other costs. Appendix Fig. S2 in [43] contains qualitative NGDF results, and App. A contains additional ablation experiments.

C. Intra-Category Generalization

To evaluate whether our method can generalize to shapes in the same object category, we train an NGDF model on 7 shapes in the "Bottle" category from ACRONYM [36]. Training samples are generated from the meshes using the same data collection procedure described in Sec. IV-C. We evaluate performance on a held-out Bottle instance, the same instance used in the previous evaluations. The intracategory model achieves 0.63 ± 0.5 execution success on 30 Bottle trials for the reaching and grasping evaluation, which is comparable with the single-object NGDF results from Table II, demonstrating intra-category generalization without loss of performance.

D. Real Robot Reaching and Grasping Evaluation

Finally, we test our method's reaching and grasping performance on a real robot system. At the start of each trial,

TABLE II: Reaching and Grasping Results

Method	Perception	Grasp Estimation	Grasp Selection	Goal	Execution Success ↑
O1 (Oracle) OMG [15] (Oracle)	Known Object Pose Known Object Pose	Known Discrete Grasps Known Discrete Grasps	Min. Distance Adaptive Cost	Fixed Variable	0.96 0.99
B1	Unknown Object Pose	Predicted Discrete Grasps [6]	Max Score	Fixed	0.37
B2	Unknown Object Pose	Predicted Discrete Grasps [6]	Min. Distance	Fixed	0.39
В3	Unknown Object Pose	Predicted Discrete Grasps [6]	Min. Distance	Variable	0.38
B4	Unknown Object Pose	Predicted Discrete Grasps [6]	Adaptive Cost	Variable	0.31
NGDF (Ours)	Unknown Object Pose	Predicted Continuous Grasps	N/A	Variable	0.61

Middle columns correspond to design decisions found in Fig. 2; color-coded methods also correspond to those shown in the same figure.

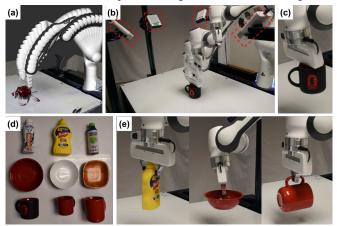


Fig. 5: Real System Evaluation. (a) Visualizing the plan and imperfect object point cloud; (b) executing the plan on hardware (cameras highlighted with red boxes); (c) lifting the object. (d) The nine objects used for testing. (e) Additional successful grasps.

an object is placed in a random stable pose. A partial point cloud of the scene is obtained from four Azure Kinect depth sensors (Fig. 5b), similar to NDF [29]. The object point cloud is segmented via plane fitting, then passed as input to NGDF models from Sec. V-A, which are trained on one instance per category in simulation. The cloud is also converted to a signed distance field to enable computing collision costs with CHOMP [11]. The optimized trajectory is executed on a Franka Panda robot with impedance control, and the trial is considered successful if the object is grasped and lifted without being dropped (Fig. 5c). 9 test objects were evaluated, 3 from each shape category (Fig. 5d). 3 grasp attempts were performed per object for a total of 27 trials. See App. C in [43] for additional details.

Our overall grasp success rate was 81%, with success per category being 7/9 Bottles, 9/9 Bowls, and 6/9 Mugs. Our system successfully grasped every object, despite many of them being outside of its training distribution in terms of size and shape. The method also demonstrated robustness to noisy perception and execution with impedance control. Failure cases were due to slight collisions between the fingers and objects, similar to what we observed in simulation.

VI. DISCUSSION

Neural implicit functions have been widely explored for 3D vision tasks such as shape reconstruction. NGDF extends this concept to grasp estimation, using 6D poses as queries on grasp manifolds. Our work differs from existing work on

3D reconstruction, not only due to the higher dimensionality of our problem, but also because of the challenge in acquiring ground truth labels. The ground truth grasp distance between an arbitrary query pose and the corresponding closest grasp is expensive to compute. Instead, we train on large-scale discrete grasp sets [36] as near-ground truth supervision. Our experiments in Sec. V-A show that NGDF is able to learn the continuous grasp manifold as the level set of the neural field from this discrete supervision.

NGDF decouples the problem of learning a grasp manifold representation from the problem of finding a good grasp pose. For the latter, we formulate the distance output of NGDF as a cost to be minimized. For the full robot motion planning regime, we jointly optimize the grasp cost with smoothness and collision costs. We outperform baselines in Sec. V-B that represent what a practitioner would implement for a reaching and grasping task. While the performance of oracle methods indicate room for improvement, our results show that joint optimization with NGDF is a promising direction for manipulation. We also demonstrate scalability with intra-category generalization results in Sec. V-C, and deploy our method on real hardware in Sec. V-D.

In terms of limitations, NGDF is trained on a gripperspecific dataset; NGDF for other grippers may require different datasets. The method also depends on upstream object segmentation. Further, the cost weights are fixed during optimization in the reach and grasp planning task; learning to adjust the weights each iteration could improve performance.

VII. CONCLUSION

We propose Neural Grasp Distance Fields (NGDF), which represent the continuous manifold of grasps as the zero-level set of a neural field. We formulate the estimated distance as a cost for a gradient-based trajectory optimizer to jointly optimize with other trajectory costs such as smoothness and collision avoidance to perform reach and grasp planning. Our results show that NGDF outperforms existing methods, while generalizing to unseen poses and unseen objects.

ACKNOWLEDGMENT

This work was supported by the US Air Force and DARPA (FA8750-18-C-0092), NSF (IIS-1849154, DGE2140739), CMU GSA/Provost Conference Funding, and the Meta AI Mentorship Program. The authors thank Kalyan Alwala and Adithya Murali for early prototyping, as well as Taosha Fan, Austin Wang, Daniel Seita, and Chuer Pan for helpful discussions and feedback.

APPENDIX

A. Ablations for Neural Grasp Distance Fields

We perform ablation experiments for our trajectory optimizer (Table S1). We compare using Adam [40] vs. a fixed step size ("No-Adam") for functional gradient descent. Unlike our method, CHOMP [11] originally uses a fixed or decaying step size, in the setting where the start and end trajectory configurations are not optimized (Sec. IV-B). In our setting, the end configuration is variable to allow optimization of the grasp pose. No-Adam converges slowly when the trajectory is far from a valid grasp pose, and overshoots when near the level set. We also evaluated using a decaying step size; while this mitigated the overshooting issue, convergence was still much slower, and the decay rate required tuning.

"No-Initial-IK" initializes the configuration at every timestep in the trajectory to the starting joint configuration, instead of using IK to initialize the trajectory as described in Sec. V-B. We observe worse performance with No-Initial-IK as the initial trajectory is farther from the desired grasp trajectory, making it harder to plan.

TABLE S1: Optimizer Ablation Results

	Grasp Execution ↑
NGDF, No-Adam	0.18
NGDF, No-Initial-IK	0.44
NGDF (Ours)	0.61

No-Adam uses CHOMP [11] with a fixed step size instead of Adam [40] optimization for the functional gradient update. No-Initial-IK initializes the trajectory so all steps in the plan start at the initial joint configuration. NGDF uses Adam and initializes the endpoint of the trajectory using inverse kinematics to achieve the best performance. 90 trials were performed as in Table II.

B. Simulation Experiment Details

1) Camera Position in Simulation: Fig. S1 shows the position of the four cameras in simulation.

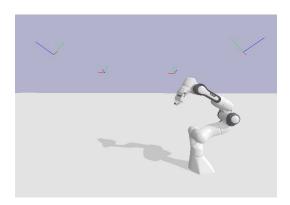


Fig. S1: Camera poses in simulation visualized as axes. The negative z axis (in blue) is the camera optical axis and points toward the robot workspace.

2) Qualitative Results: Fig. S2 visualizes successful grasp trajectories in simulation for the reaching and grasping task from Sec. V-A.



Fig. S2: Successful grasp trajectories (left-to-right) planned by our method for the bowl (top) and mug (bottom).

C. Real System Experiment Details

This section provides system implementation details for our real world experiments in Sec. V-D. Our system consists of a 7-DOF Franka Panda robot and four Azure Kinect cameras (Fig. 5b), a similar setup to NDF [29].

- 1) Calibration: The Azure Kinect cameras were extrinsically calibrated using ColoredICP [44]. For camera intrinsics, the factory calibration was used. Robot-camera extrinsic calibration was performed using Tsai-Lenz [45]. The calibrated cameras produce a combined scene point cloud in the robot base frame.
- 2) Point Cloud Processing: To segment the object point cloud from the scene point cloud, we fit a table plane using RANSAC and remove points belonging to the plane. Outlier removal and DBScan [46] are used to refine the object point cloud. Our planner requires a signed distance field (SDF) of the object for collision avoidance, so we construct a mesh from the object point cloud, then compute the SDF from the mesh using the tools provided in Wang et al. [15].

Even with four cameras, careful calibration, and point cloud processing, we recover partial point clouds with inaccuracies and noise (see Fig. S3). Despite these deficiencies, our method achieves a high success rate on real objects in various configurations (Sec. V-D), demonstrating robustness to perceptual errors.

3) Control: The output of the planner is a joint angle trajectory consisting of 30 timesteps. In order to execute the trajectory on the Franka Panda, the total duration of trajectory execution is set to 5 seconds, and the trajectory is interpolated using cubic spline interpolation to provide joint angles at 1 Hz. Impedance control [47] is then used to execute the high-frequency trajectory.

D. Experimental Setup

Several of the test objects could be grasped via multiple stable pose configurations. For example, a mug can be grasped while upright, on its side, or upside down. For objects with multiple graspable pose configurations, the configuration was randomly sampled for each trial. The 9 test objects (see Fig. 5d) had graspable stable pose configurations shown in Table S2.

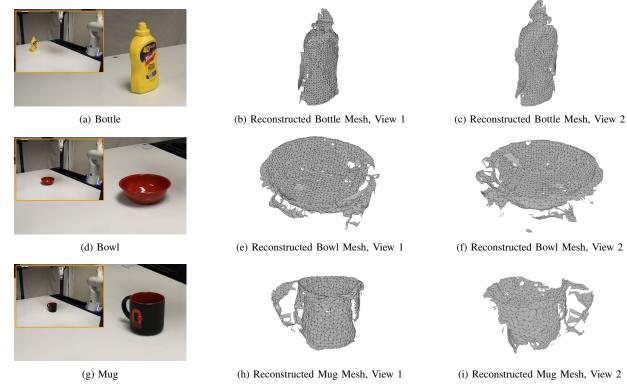


Fig. S3: Meshes reconstructed during system evaluations. The first column shows the placement of the objects in each trial, along with a close-up image of the object itself. The second column shows the meshes reconstructed from the four depth cameras according to the procedure in Sec. C.2, posed roughly as they appear in the first column. The third column shows the back of each mesh. Note that these meshes and images are magnified for visual clarity and are not consistently scaled. Even with outlier removal and other mesh processing techniques, we observe inaccuracies in the reconstruction; however, our method is robust to these inaccuracies as demonstrated by our results in Sec. V-D.

TABLE S2: Real Object Pose Configurations

	Sampled Pose Configurations
Bottle 1 (Protein Drink)	Upright, Sideways
Bottle 2 (Mustard Bottle)	Upright
Bottle 3 (Coconut Water)	Upright, Sideways
Bowl 1 (YCB Bowl)	Upright
Bowl 2 (White Bowl)	Upright
Bowl 3 (Square Bowl)	Upright
Mug 1 (Black Mug)	Upright, Sideways, Upside Down
Mug 2 (YCB Mug)	Upright, Sideways, Upside Down
Mug 3 (Large Mug)	Upright, Sideways, Upside Down

The sampled pose configurations for objects in the real system evaluation. Objects are numbered left to right according to Fig. 5d. Some stable poses did not permit grasping and were omitted; for example, the mustard bottle cannot be grasped lying sideways as it is too wide.

REFERENCES

- B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in ECCV, 2020.
- [2] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, "Deepsdf: Learning continuous signed distance functions for shape representation," in *Proceedings of the IEEE/CVF Conference on Com*puter Vision and Pattern Recognition, 2019, pp. 165–174.
- [3] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, "Occupancy networks: Learning 3d reconstruction in function space,"

- in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 4460–4470.
- [4] J. Chibane, A. Mir, and G. Pons-Moll, "Neural unsigned distance fields for implicit function learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, December 2020.
- [5] J. Chibane, T. Alldieck, and G. Pons-Moll, "Implicit functions in feature space for 3d shape reconstruction and completion," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2020.
- [6] M. Sundermeyer, A. Mousavian, R. Triebel, and F. Dieter, "Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes," *IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
- [7] A. Mousavian, C. Eppner, and D. Fox, "6-dof graspnet: Variational grasp generation for object manipulation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2901–2910.
- [8] P. Ni, W. Zhang, X. Zhu, and Q. Cao, "Pointnet++ grasping: Learning an end-to-end spatial grasp generation algorithm from sparse point clouds," in 2020 IEEE International Conference on Robotics and Automation (ICRA), 2020, pp. 3619–3625.
- [9] H. Liang, X. Ma, S. Li, M. Görner, S. Tang, B. Fang, F. Sun, and J. Zhang, "Pointnetgpd: Detecting grasp configurations from point sets," in 2019 International Conference on Robotics and Automation (ICRA). IEEE, 2019, pp. 3629–3635.
- [10] H.-S. Fang, C. Wang, M. Gou, and C. Lu, "Graspnet-Ibillion: A large-scale benchmark for general object grasping," in *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition(CVPR), 2020, pp. 11 444–11 453.
- [11] M. Zucker, N. Ratliff, A. D. Dragan, M. Pivtoraiko, M. Klingensmith, C. M. Dellin, J. A. Bagnell, and S. S. Srinivasa, "Chomp: Covariant hamiltonian optimization for motion planning," *The International Journal of Robotics Research*, vol. 32, no. 9-10, pp. 1164–1193, 2013.
- [12] J. Kuffner and S. LaValle, "Rrt-connect: An efficient approach to single-query path planning," in *Proceedings 2000 ICRA*. Millennium

- Conference. IEEE International Conference on Robotics and Automation. Symposia Proceedings (Cat. No.00CH37065), vol. 2, 2000, pp. 995–1001 vol.2.
- [13] L. Kavraki, P. Svestka, J.-C. Latombe, and M. Overmars, "Probabilistic roadmaps for path planning in high-dimensional configuration spaces," *IEEE Transactions on Robotics and Automation*, vol. 12, no. 4, pp. 566–580, 1996.
- [14] M. Mukadam, J. Dong, X. Yan, F. Dellaert, and B. Boots, "Continuoustime gaussian process motion planning via probabilistic inference," *The International Journal of Robotics Research*, vol. 37, no. 11, pp. 1319–1340, 2018.
- [15] L. Wang, Y. Xiang, and D. Fox, "Manipulation trajectory optimization with online grasp synthesis and selection," in *Robotics: Science and Systems (RSS)*, 2020.
- [16] D. Berenson, S. Srinivasa, and J. Kuffner, "Task space regions: A framework for pose-constrained manipulation planning," *The Interna*tional Journal of Robotics Research, vol. 30, no. 12, pp. 1435–1460, 2011.
- [17] J. Bohg, A. Morales, T. Asfour, and D. Kragic, "Data-driven grasp synthesis—a survey," *IEEE Transactions on Robotics*, vol. 30, no. 2, pp. 289–309, 2014.
- [18] K. Kleeberger, R. Bormann, W. Kraus, and M. F. Huber, "A survey on learning-based robotic grasping," *Current Robotics Reports*, vol. 1, no. 4, pp. 239–249, 2020.
- [19] J. Ichnowski, M. Danielczuk, J. Xu, V. Satish, and K. Goldberg, "Gomp: Grasp-optimized motion planning for bin picking," in 2020 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2020, pp. 5270–5277.
- [20] A. D. Dragan, N. D. Ratliff, and S. S. Srinivasa, "Manipulation planning with goal sets using constrained trajectory optimization," in 2011 IEEE International Conference on Robotics and Automation, 2011, pp. 4582–4588.
- [21] L. Wang, X. Meng, Y. Xiang, and D. Fox, "Hierarchical policies for cluttered-scene grasping with latent plans," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 2883–2890, 2022.
- [22] S. Song, A. Zeng, J. Lee, and T. Funkhouser, "Grasping in the wild: Learning 6dof closed-loop grasping from low-cost demonstrations," *IEEE Robotics and Automation Letters*, vol. 5, no. 3, pp. 4978–4985, 2020.
- [23] W. Yang, C. Paxton, A. Mousavian, Y.-W. Chao, M. Cakmak, and D. Fox, "Reactive human-to-robot handovers of arbitrary objects," *IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
- [24] K. Karunratanakul, J. Yang, Y. Zhang, M. J. Black, K. Muandet, and S. Tang, "Grasping field: Learning implicit representations for human grasps," in 2020 International Conference on 3D Vision (3DV). IEEE, 2020, pp. 333–344.
- [25] L. Zhu, A. Mousavian, Y. Xiang, H. Mazhar, J. van Eenbergen, S. Debnath, and D. Fox, "Rgb-d local implicit function for depth completion of transparent objects," in *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, 2021, pp. 4649–4658
- [26] J. Ichnowski*, Y. Avigal*, J. Kerr, and K. Goldberg, "Dex-NeRF: Using a neural radiance field to grasp transparent objects," in *Conference on Robot Learning (CoRL)*, 2020.
- [27] Y. Wi, P. Florence, A. Zeng, and N. Fazeli, "Virdo: Visio-tactile implicit representations of deformable objects," 2022.
- [28] P. Florence, C. Lynch, A. Zeng, O. A. Ramirez, A. Wahid, L. Downs, A. Wong, J. Lee, I. Mordatch, and J. Tompson, "Implicit behavioral cloning," in *Conference on Robot Learning*. PMLR, 2022, pp. 158–
- [29] A. Simeonov, Y. Du, A. Tagliasacchi, J. B. Tenenbaum, A. Ro-driguez, P. Agrawal, and V. Sitzmann, "Neural descriptor fields: Se(3)-

- equivariant object representations for manipulation," arXiv preprint arXiv:2112.05124, 2021.
- [30] Y. Li, S. Li, V. Sitzmann, P. Agrawal, and A. Torralba, "3d neural scene representations for visuomotor control," in *Conference on Robot Learning*. PMLR, 2022, pp. 112–123.
- [31] D. Driess, J.-S. Ha, M. Toussaint, and R. Tedrake, "Learning models as functionals of signed-distance fields for manipulation planning," in *Conference on Robot Learning*. PMLR, 2022, pp. 245–255.
- [32] Z. Jiang, Y. Zhu, M. Svetlik, K. Fang, and Y. Zhu, "Synergies between affordance and geometry: 6-dof grasp detection via implicit representations," 2021.
- [33] Y.-H. Wu, J. Wang, and X. Wang, "Learning generalizable dexterous manipulation from human grasp affordance," Conference on Robot Learning, 2022.
- [34] N. Khargonkar, N. Song, Z. Xu, B. Prabhakaran, and Y. Xiang, "Neuralgrasps: Learning implicit representations for grasps of multiple robotic hands," *Conference on Robot Learning*, 2022.
- [35] J. Urain, N. Funk, G. Chalvatzaki, and J. Peters, "Se (3)-diffusionfields: Learning cost functions for joint grasp and motion optimization through diffusion," arXiv preprint arXiv:2209.03855, 2022.
- [36] C. Eppner, A. Mousavian, and D. Fox, "ACRONYM: A large-scale grasp dataset based on simulation," in *IEEE International Conference* on Robotics and Automation (ICRA), 2021.
- [37] M. Macklin, M. Müller, N. Chentanez, and T.-Y. Kim, "Unified particle physics for real-time applications," ACM Trans. Graph., vol. 33, no. 4, jul 2014. [Online]. Available: https://doi.org/10.1145/ 2601097.2601152
- [38] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," http://pybullet.org, 2016–2021
- [39] C. Deng, O. Litany, Y. Duan, A. Poulenard, A. Tagliasacchi, and L. J. Guibas, "Vector neurons: A general framework for so (3)-equivariant networks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12200–12209.
- [40] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," CoRR, vol. abs/1412.6980, 2015.
- [41] L. Pineda, T. Fan, M. Monge, S. Venkataraman, P. Sodhi, R. Chen, J. Ortiz, D. DeTone, A. Wang, S. Anderson et al., "Theseus: A library for differentiable nonlinear optimization," arXiv preprint arXiv:2207.09442, 2022.
- [42] G. Sutanto, A. Wang, Y. Lin, M. Mukadam, G. Sukhatme, A. Rai, and F. Meier, "Encoding physical constraints in differentiable newton-euler algorithm," ser. Proceedings of Machine Learning Research, A. M. Bayen, A. Jadbabaie, G. Pappas, P. A. Parrilo, B. Recht, C. Tomlin, and M. Zeilinger, Eds., vol. 120. The Cloud: PMLR, 10–11 Jun 2020, pp. 804–813. [Online]. Available: http://proceedings.mlr.press/v120/sutanto20a.html
- [43] T. Weng, D. Held, F. Meier, and M. Mukadam, "Neural grasp distance fields for robot manipulation," arXiv preprint arXiv:2211.02647, 2022.
- [44] J. Park, Q.-Y. Zhou, and V. Koltun, "Colored point cloud registration revisited," in *Proceedings of the IEEE international conference on* computer vision, 2017, pp. 143–152.
- [45] R. Y. Tsai, R. K. Lenz et al., "A new technique for fully autonomous and efficient 3 d robotics hand/eye calibration," *IEEE Transactions on* robotics and automation, vol. 5, no. 3, pp. 345–358, 1989.
- [46] M. Ester, H.-P. Kriegel, J. Sander, X. Xu et al., "A density-based algorithm for discovering clusters in large spatial databases with noise." in kdd, vol. 96, no. 34, 1996, pp. 226–231.
- [47] K. Zhang, M. Sharma, J. Liang, and O. Kroemer, "A modular robotic arm control stack for research: Franka-interface and frankapy," arXiv preprint arXiv:2011.02398, 2020.