Neural Motion Fields: Encoding Grasp Trajectories as Implicit Value Functions

Yun-Chun Chen^{1,2,*} Adithyavairavan Murali^{1,*} Balakumar Sundaralingam^{1,*}
Wei Yang¹ Animesh Garg^{1,2} Dieter Fox^{1,3}

¹NVIDIA ²University of Toronto ³University of Washington *Equal contribution

Abstract—The pipeline of current robotic pick-and-place methods typically consists of several stages: grasp pose detection, finding inverse kinematic solutions for the detected poses, planning a collision-free trajectory, and then executing the open-loop trajectory to the grasp pose with a low-level tracking controller. While these grasping methods have shown good performance on grasping static objects on a table-top, the problem of grasping dynamic objects in constrained environments remains an open problem. We present Neural Motion Fields, a novel object representation which encodes both object point clouds and the relative task trajectories as an implicit value function parameterized by a neural network. This object-centric representation models a continuous distribution over the SE(3) space and allows us to perform grasping reactively by leveraging sampling-based MPC to optimize this value function.

I. INTRODUCTION

Current robotic grasping approaches typically decompose the task of grasping into several sub-components: detecting grasp poses on point clouds [13, 14], finding inverse kinematics solutions at these poses, solving collision-free trajectories to pre-grasp standoff poses and finally executing open-loop trajectories from standoff poses to the grasp poses [8, 9]. By inferring a finite discrete number of grasp poses, such an approach neglects the insight that object grasp affordances are a continuous manifold. While this approach has yielded tremendous progress in bin-picking and grasping unknown objects on a table-top, reactive grasping of unknown objects in constrained environments remains an open problem.

Implicit neural representations [6, 7, 10] have emerged as a new paradigm for applications in rendering, view synthesis, and shape reconstruction. Compared to traditional explicit representations (e.g., point clouds and meshes), implicit neural representations can represent continuous signals at arbitrarily high resolutions. Motivated by implicit neural representations, we propose to learn a value function that encodes robotic task trajectories using a neural network. Our key insight is to map each gripper pose in SE(3) to its trajectory path length as shown in Figure 1. To train the model, we generate synthetic data of grasping process by using a prior grasp dataset [8] and planning trajectories with a RRT [5] motion planner. Once the model training is done, we cast the learned value function as a cost and leverage the Model Predictive Control (MPC) [1] algorithm to query gripper poses in SE(3) with cost minimization. This allows us to generate a kinematically feasible trajectory that the robot can execute to reach a grasp on the object.

We benchmark our method on the grasping task and report the success rate. In addition, we evaluate our model in two

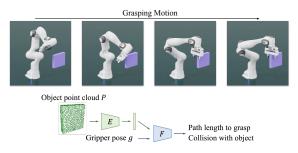


Fig. 1: We present Neural Motion Fields that learns a value function that can be queried to generate grasping motion. We use separate network weights for predicting path lengths and collisions.

settings: static object poses and dynamic object poses, and provide ablation studies under various settings. In summary, our contributions include

- 1) We propose Neural Motion Fields, a novel formulation of the grasp motion generation problem in SE(3) as a continuous implicit representation.
- 2) We show that this learned object-centric representations allows reactive grasp manipulation using MPC [1].

II. PROBLEM STATEMENT

Value function learning for grasping. We are interested in learning a model that can be used to plan a kinematically feasible trajectory for the robot to execute to grasp an object. Specifically, we cast this task as a value function learning problem. We assume that we are given a segmented object point cloud $P \in \mathbb{R}^{N \times 3}$, where N is the number of points in a point cloud, and a gripper pose $g \in SE(3)$. The value function V(g,P) describes how far the gripper pose g is from a grasp on the object. We use the path length of a gripper pose to represent the value function.

Gripper pose path length. As shown in Figure 2a, given a trajectory $\{g_i\}_{i=0}^t$, where g_0 denotes the end pose (grasp pose) and g_t denotes the start pose, the path length of the start pose g_t (denoted as $V(g_t)$) is defined as the cumulative sum of the average distance between two adjacent gripper poses [16], which can be expressed by

$$V(g_t) = \sum_{i=0}^{t-1} \frac{1}{m} \sum_{x \in M} \| (R_i x + T_i) - (R_{i+1} x + T_{i+1}) \|, \quad (1)$$

where R_i is the rotation of the gripper pose g_i , T_i is the translation of the gripper pose g_i , m is the number of keypoints of the gripper, and M is the set of keypoints of the gripper.

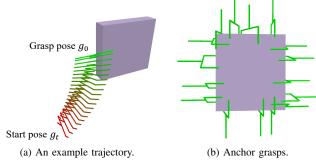


Fig. 2: (a) An example trajectory planned by RRT. Red represents longer path length. Green represents shorter path length. (b) Visualization of the selected anchor grasps.

III. NEURAL MOTION FIELDS

A. Learning from Grasp Trajectories

Given an object point cloud P and a gripper pose g, our goal is to learn a model that approximates the value function V(g,P). We propose Neural Motion Fields, which consists of two modules: a path length module and a collision module.

Path length prediction. As shown in Figure 1, the path length module first takes as input the object point cloud P and uses the point cloud encoder $E_{\text{path-length}}$ to encode a feature embedding $f_{\text{path-length}} = E_{\text{path-length}}(P) \in \mathbb{R}^d$, where d is the dimension of the feature $f_{\text{path-length}}$. Then, the point cloud feature $f_{\text{path-length}}$ and the gripper pose g are concatenated and passed to the path length prediction network $F_{\text{path-length}}$ to predict the path length $V_{\text{pred}}(g,P)$ for the gripper pose g.\frac{1}{}

To train the path length module, we adopt an ℓ_1 loss function, which is defined as

$$\mathcal{L}_{\text{path-length}} = ||V_{\text{pred}}(g, P) - V_{\text{gt}}(g, P)||_1, \tag{2}$$

where $V_{\text{pred}}(g, P)$ denotes the predicted path length of the gripper pose g and $V_{\text{gt}}(g, P)$ denotes the ground truth.

We visualize the learned value function using a cost map visualization as shown in Figure 3. We show two cost maps. In each cost map, we select a grasp pose. We keep the orientation and vary the x and y positions of the grasp pose to compose poses. We then query the path lengths of the composed poses using the learned model. Each input pose is represented by a point in \mathbb{R}^3 and is colored by its predicted path length (red means longer path length, while green means shorter).

Collision prediction. Having the path length module alone is insufficient as the model does not explicitly penalize collisions between the gripper and the object of interest. To address this issue, we develop a collision module as shown in Figure 1 (same input as the path length module, but mapped to a different output using a different set of network weights).

Given an object point cloud P and a gripper pose g, the collision module first uses the point cloud encoder $E_{\text{collision}}$ to encode the point cloud feature $f_{\text{collision}} = E_{\text{collision}}(P) \in \mathbb{R}^k$,

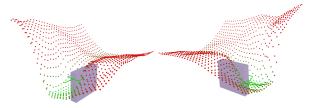


Fig. 3: We visualize two cost maps of a box object. In each cost map, we select a grasp pose. We keep the orientation, vary the x and y positions, and query the path length of the composed gripper pose using our model. Each gripper pose is represented by a point in 3D. Red represents longer path length (higher in z). Green represents shorter path length (lower in z).

where k is the dimension of the feature $f_{\rm collision}$. Then, the point cloud feature $f_{\rm collision}$ and the gripper pose g are concatenated and passed to the collision prediction network $F_{\rm collision}$ to predict the probability $p_{\rm pred}(g,P)$ of the gripper pose g being in collision with the object.

To train the collision model, we adopt a standard binary cross-entropy loss function, which is defined as

$$\mathcal{L}_{\text{collision}} = p_{\text{gt}}(g, P) \log p_{\text{pred}}(g, P) + (1 - p_{\text{gt}}(g, P)) \log(1 - p_{\text{pred}}(g, P)),$$
(3)

where $p_{\text{pred}}(g, P)$ is the predicted collision probability and $p_{\text{gt}}(g, P)$ is the ground truth.

B. Generating Grasp Motion

Given the path length value function represented by the path length module and the collision value function represented by the collision module, we formulate the grasp cost C_{grasp} as

$$C_{\text{grasp}}(g_t, P) = (1 - V(g_t, P)) + C(g_t, P),$$
 (4)

where $V(g_t, P)$ is the predicted path length for the gripper pose g_t , $C(g_t, P)$ is the collision cost of the gripper pose g_t computed by thresholding $p(g_t, P) \ge \tau$, τ is a hyperparameter, and P is the object point cloud. In our work, we set $\tau = 0.25$.

We then optimize the grasp cost C_{grasp} along with the cost C_{storm} to ensure smooth collision-free motions using STORM [1], which is a GPU-based MPC framework:

$$\min_{\ddot{x}_{t \in [0,H]}} \quad \mathcal{C}_{\text{storm}}(q) + \mathcal{C}_{\text{grasp}} \tag{5}$$

Additional details on C_{storm} is available in [1].

IV. EXPERIMENTS

A. Experimental Setup

Dataset. We experiment with four box objects from the dataset provided by [8] and one bowl object from the ACRONYM dataset [3]. We first subsample a set of 16 grasp poses around the object using farthest point sampling as shown in Figure 2b. We then apply a 180 degree rotation along the z-axis of the gripper pose to get the flipped counterparts. The selected 32 grasps are then used as the goal poses for trajectory data collection. We randomly sample a gripper pose with respect to one of the goal poses and use the RRT [5] planner from OMPL [12] to plan a trajectory between the two poses. We note that our data generation pipeline is agnostic to the choice of the motion planning algorithm and other planners are also

 $^{^1}$ The gripper pose g is first converted to a vector in \mathbb{R}^9 which is the contatenatation of the 6D rotation representation [17] of g and the translation vector in \mathbb{R}^3 of g. The 9D vector will then be concatenated with the point cloud feature $f_{\text{path-length}}$ for path length prediction.

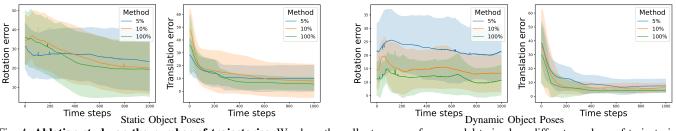


Fig. 4: Ablation study on the number of trajectories. We show the rollout curves of our model trained on differnt numbers of trajectories.

TABLE I: Grasp success rate.

| Method | Bowl | Box A | Box B | Box C | Box D |
|----------------|------|-------------|-------|------------|------------|
| Oracle Ours | | 100% 80% | | 30% 40% | 50% 30% |

applicable. We use inverse kinematics to solve for the joint angles for each waypoint returned by RRT and interpolate between two adjacent joint angles to obtain denser waypoints for each trajectory. For each object, we collect one million trajectories. We further filter out trajectories where the start gripper pose has a path length greater than a threshold $\phi=1.0$.

Implementation details. We implement our model using PyTorch [11]. We use the ADAM [4] optimizer for model training. We use DGCNN [15] to be our point cloud encoder. The path length prediction network $F_{\text{path-length}}$ and the collision prediction network $F_{\text{collision}}$ both consist of 20 fully connected layers. The learning rate is set to 2×10^{-3} with a weight decay of 1×10^{-6} . We train our model using eight NVIDIA V100 GPUs with 32GB memory each. The batch size is set to 32. The number of points in a point cloud is set to 1,024. The dimensions for the point cloud features $f_{\text{path-length}}$ and $f_{\text{collision}}$ are both 512. The training time for both the path length module and the collision module is around 7 days.

B. Evaluation on Grasping

Setting. For each object, there are 10 test cases. In each test case, we initialize the object with a stable pose on the tabletop. In each test case, we optimize over the learned path length module and the collision module to find the minimum point by leveraging sampling based optimization [1]. The optimized gripper pose is then used as the grasp pose. We use STORM pose reaching [1] to reach the grasp pose. We set the time limit for each test case to 30 seconds.

Metric. We use the grasp success rate to evaluate performance. **Results.** Table I reports the grasp success rate of the 5 objects. Our model performs well on Box A, but achieves inferior performance on all other objects. The inferior performance is due to the collision between the finger tips of the gripper and the object.

C. Ablation Study

Setting. We conduct two ablation studies: 1) static object pose reaching and 2) dynamic object pose reaching. In the static object pose reaching setting, we randomly sample an object pose at the beginning of each episode and place the object at the sampled pose. The object pose is kept static throughout the episode. In the dynamic object pose reaching setting, we

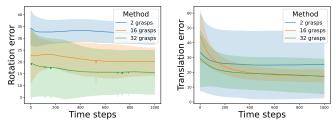


Fig. 5: **Ablation study on the number of anchor grasps.** We show the rollout curves of our model trained on trajectories collected by using differnt numbers of anchor grasps.

randomly sample an object pose and a velocity vector at the beginning of each episode and place the object at the sampled pose while the object is moving at the speed specified by the velocity vector. In both settings, there are 50 episodes. Each episode has 1,000 time steps. We run MPC with our model to optimize the grasp cost in Equation (5) and execute the robot.

Metric. We follow [2] and compute the rotation error and the translation error between the current gripper pose and the closest grasp pose at each time step. The closest grasp pose is the grasp pose that has the minimum path length from the current gripper pose.

Ablation study on the number of trajectories. Figure 4 shows the rollout curves of our model trained on different percentages of the collected dataset (i.e., 5%, 10% and 100%). We observe that in both settings, training on the entire dataset (i.e., 100%) achieves the best in both the rotation error and the translation error.

Abaltion study on the number of grasp poses. Figure 5 shows the rollout curves of our model trained on trajectories collected by using different numbers of anchor grasps (i.e., 2, 16 and 32 grasps) in the dynamic object pose reaching setting. We observe that the model trained on trajectories generated from 32 anchor grasps achieves the best performance and converges faster in translation error than all other models.

V. CONCLUSIONS AND FUTURE WORK

We propose Neural Motion Fields, a novel object representation which encodes both object point clouds and the relative task trajectories as an implicit value function parameterized by a neural network. This object-centric representation models a continuous distribution over the SE(3) space and allows us to perform grasping reactively by leveraging sampling-based MPC to optimize this value function. Through experimental evaluations, we show that by training on more numbers of anchor grasps and larger scale datasets results in superior performance. In future work, we plan to train a single model on more objects.

REFERENCES

- [1] Mohak Bhardwaj, Balakumar Sundaralingam, Arsalan Mousavian, Nathan D Ratliff, Dieter Fox, Fabio Ramos, and Byron Boots. STORM: An Integrated Framework for Fast Joint-Space Model-Predictive Control for Reactive Manipulation. In *CoRL*, 2022. 1, 2, 3
- [2] Silvia Cruciani, Balakumar Sundaralingam, Kaiyu Hang, Vikash Kumar, Tucker Hermans, and Danica Kragic. Benchmarking in-hand manipulation. *IEEE Robotics and Automation Letters*, 2020. 3
- [3] Clemens Eppner, Arsalan Mousavian, and Dieter Fox. ACRONYM: A large-scale grasp dataset based on simulation. In *ICRA*, 2021. 2
- [4] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014. 3
- [5] Steven M LaValle, James J Kuffner, BR Donald, et al. Rapidly-exploring random trees: Progress and prospects. *Algorithmic and computational robotics: new directions*, 2001. 1, 2
- [6] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In CVPR, 2019. 1
- [7] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In ECCV, 2020. 1
- [8] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof graspnet: Variational grasp generation for object manipulation. In *ICCV*, 2019. 1, 2
- [9] Adithyavairavan Murali, Arsalan Mousavian, Clemens Eppner, Chris Paxton, and Dieter Fox. 6-dof grasping for target-driven object manipulation in clutter. In *ICRA*, 2020. 1
- [10] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In CVPR, 2019. 1
- [11] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In NeurIPS, 2019. 3
- [12] Ioan A Sucan, Mark Moll, and Lydia E Kavraki. The open motion planning library. *IEEE Robotics & Automation Magazine*, 2012. 2
- [13] Martin Sundermeyer, Arsalan Mousavian, Rudolph Triebel, and Dieter Fox. Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes. In *ICRA*, 2021. 1
- [14] Lirui Wang, Yu Xiang, Wei Yang, Arsalan Mousavian, and Dieter Fox. Goal-Auxiliary Actor-Critic for 6D Robotic Grasping with Point Clouds. In *CoRL*, 2021. 1

- [15] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. ACM Transactions on Graphics, 2019. 3
- [16] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. In RSS, 2018. 1
- [17] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *CVPR*, 2019. 2