

DexNoMa: Learning Geometry-Aware Nonprehensile Dexterous Manipulation

Anonymous Author(s)

Affiliation

Address

email

Abstract: Nonprehensile manipulation, such as pushing and pulling, enables robots to move, align, or reposition objects that may be difficult to grasp due to their geometry, size, or relationship to the robot or the environment. Much of the existing work in nonprehensile manipulation relies on parallel-jaw grippers or tools such as rods and spatulas. Multi-fingered dexterous hands offer richer contact modes and versatility for handling diverse objects to provide stable support over the objects, which compensates for the difficulty of modeling the dynamics of nonprehensile manipulation. We propose **Dexterous Nonprehensile Manipulation** (DexNoMa), a method for nonprehensile manipulation which frames the problem as synthesizing and learning pre-contact dexterous hand poses that lead to effective pushing and pulling. We generate diverse hand poses via contact-guided sampling, filter them using physics simulation, and train a diffusion model conditioned on object geometry to predict viable poses. At test time, we sample hand poses and use standard motion planning tools to select and execute pushing and pulling actions. We perform 840 real-world experiments with an Allegro Hand, comparing our method to baselines. The results indicate that DexNoMa offers a scalable route for training dexterous nonprehensile manipulation policies. Our pre-trained models and dataset, including 1.3 million hand poses across 2.3k objects, will be open-source to facilitate further research. Supplementary material is available here: dexnoma.github.io.

Keywords: Nonprehensile manipulation, dexterous hand

1 Introduction

Nonprehensile actions are fundamental to how humans and robots interact with the physical world [1, 2, 3, 4]. These actions permit the manipulation of objects that may be too large, heavy, or geometrically complex to grasp directly. While there has been tremendous progress in nonprehensile robot manipulation [5, 6, 7, 8, 9], most work uses simple end-effectors such as parallel-jaw grippers, rods [10, 11], or spatulas [12]. In contrast, multi-fingered hands with high degrees-of-freedom (DOF) such as the Allegro Hand or LEAP Hand [13] enable contact patterns that can be especially useful for stabilizing complex, awkward, or top-heavy objects, or for coordinating contact across multiple objects. However, despite their promise and recent progress [14], leveraging high-DOF hands for nonprehensile manipulation remains relatively underexplored due to the challenges of modeling hand-object relationships and planning feasible contact-rich motions.

In this paper, we study pushing and pulling objects using the 4-finger, 16-DOF Allegro Hand. Our insight is to recast this problem into one of synthesizing effective pre-contact hand poses, an approach inspired by recent success in generating large-scale datasets for dexterous *grasping* [15, 16, 17, 18, 19, 20]. We propose a scalable pipeline for generating hand poses for pushing and pulling objects. This involves contact-guided optimization and validation via GPU-accelerated physics simulation with IsaacGym [21]. These filtered hand poses are then used to train a generative diffusion policy conditioned on object geometry, represented using basis point sets [22].

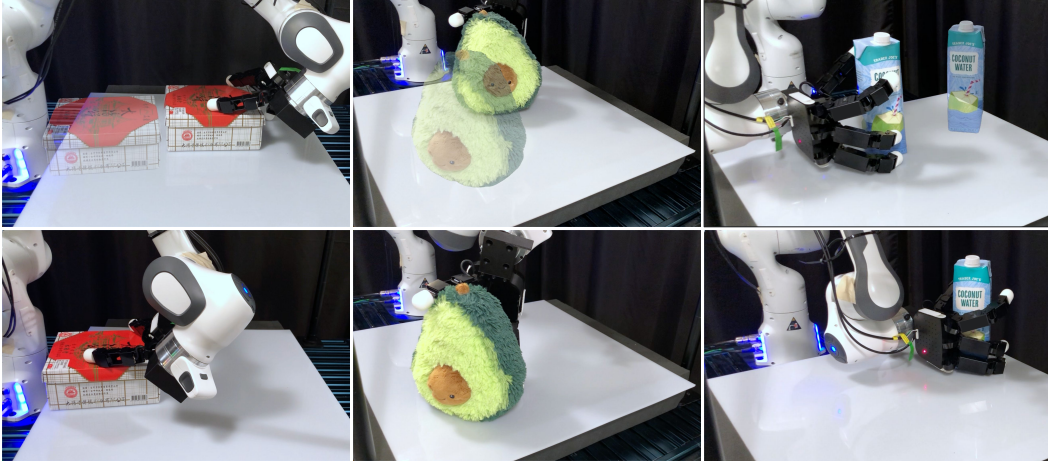


Figure 1: Three examples (one per column) of nonprehensile manipulation using DexNoMa with a 4-finger, 16-DOF Allegro Hand. The top row shows the starting object configuration with its goal rendered as a transparent overlay, while the bottom row shows the result after the robot’s motion. DexNoMa synthesizes diverse hand poses conditioned on object geometry, handling flat (left), volumetric (middle), and tall (right) objects.

At test time, we use visual data to reconstruct an object mesh in physics simulation. The trained diffusion policy uses this mesh to generate diverse hand poses for pushing or pulling. We then validate the resulting hand poses in simulation, and execute the best-performing action in the real world. We call this pipeline **Dexterous Nonprehensile Manipulation (DexNoMa)**. Figure 1 shows several real-world examples where the hand pose differs depending on object geometry. Overall, our experimental results across diverse common and 3D-printed objects demonstrate that DexNoMa is a promising approach for generalizable object pushing and pulling. It outperforms alternative methods such as querying the nearest hand pose in our data or using a fixed spatula-like hand pose, highlighting the need for a diffusion model to generate diverse hand poses.

To summarize, the contributions of this paper include:

- A scalable pipeline for generating and filtering dexterous hand poses for pushing and pulling.
- A diffusion model for geometry-conditioned hand pose prediction for nonprehensile manipulation.
- A motion planning framework to execute these poses for nonprehensile manipulation in the real world, with results across **840 trials** showing that DexNoMa outperforms alternative methods.
- A dataset of 1.3 million hand poses for pushing and pulling across 2.3k objects with corresponding canonical point cloud observations.

2 Related Work

Nonprehensile Robot Manipulation. Classical nonprehensile manipulation includes pushing, sliding, rolling, and tilting, and has a long history in robotics [1, 2, 3, 4]. Planning methods for nonprehensile manipulation often assume access to object models or priors [23, 24, 25]. Another recent planning-based method explores nonprehensile interaction with high-DOF hands in simulation by analyzing contact reasoning and wrench closure [26]. In contrast, our work targets real-world pushing and pulling using a high-DOF hand applied to diverse and geometrically complex objects. Recent learning-based methods have extended nonprehensile manipulation beyond classical planning, including extrinsic dexterity systems [5, 27] and those based on predicting object dynamics such as HACMan [6, 7], CORN [8], and DyWA [9]. Other works approach pushing as a precursor to grasping, often in planar settings with simple parallel-jaw grippers for multi-object manipulation [28, 29], or use bimanual systems for nonprehensile tasks using multi-link tools [30]. None of these works study learning for single-hand pushing and pulling with dexterous hands. Furthermore, many prior benchmarks focus on pushing single flat objects on a surface, such as a T-shape object [11], or use spatulas to move small cubes and granular media [12, 10]. Our work directly targets larger and more complex objects, including those that might topple or require coordinated multi-surface contact.

Dexterous Grasping Synthesis and Datasets. A substantial body of research focuses on generating and evaluating grasp poses for multi-fingered hands. Pioneering efforts such as Liu et al. [31] create a dataset of 6.9K grasps using the GraspIt! [32] software tool, while Jiang et al. [33] synthesize human hand poses by using a conditional Variational Autoencoder [34]. More recent efforts significantly scale grasp generation with tools such as differentiable contact simulation [35, 36] or optimization over an energy function based on Differentiable Force Closure (DFC) [37]. Our work falls in the latter category, which has facilitated the generation of diverse grasping datasets such as DexGraspNet [16] with 1.32M grasps followed by DexGraspNet 2.0 [17] with 427M grasps. These pipelines generate hand poses by optimization over an energy function, filter them using physics simulators, train generative diffusion models for grasp synthesis, and typically include some fine-tuning or evaluation modules [38, 15]. While our pipeline also uses energy-based pose optimization and filtering, our focus is on generating hand poses for nonprehensile manipulation.

Learning-Based Dexterous Manipulation. Learning-based approaches for robotic grasping and manipulation have rapidly expanded in recent years [39, 40]. While some recent work emphasizes fine-grained bimanual manipulation using parallel-jaw grippers [41, 42], our focus is on learning single-arm manipulation with high-DOF dexterous hands such as the LEAP [13], Allegro, and Shadow hands. These hands have been applied to a variety of tasks, such as in-hand object rotation [43, 44, 45, 46, 47], object singulation [48, 49], multi-object manipulation [50, 20, 51, 29], and bimanual systems [52, 53]. While showing the versatility of dexterous hardware, these works focus on largely prehensile interactions. Prior learning-based systems with high-DOF hands for non-prehensile behaviors demonstrate tasks such as rolling objects or picking up plates as examples of learning from 3D data [54] or human videos [55]. Recently, Chen et al. [56] synthesize task-oriented dexterous hand poses for certain nonprehensile tasks such as pulling drawers. However, none of these methods directly study pushing or pulling as their primary manipulation mode.

3 Problem Statement and Assumptions

We study nonprehensile object manipulation on a flat surface using a single-arm robot with a high-DOF multi-finger dexterous hand (e.g., the Allegro Hand). By “nonprehensile,” we specifically refer to *pushing* or *pulling* in this paper. We assume that there exists one object O on the surface with configuration $S_{\text{obj}} \in SE(3)$, and that the surface’s friction properties facilitate object pushing. We use P to indicate the object’s point cloud sampled from its surface. Let \mathcal{H} be the space of possible nonprehensile hand poses, where $H \in \mathcal{H}$ is defined as $H = (\theta, T)$. Here, $\theta \in \mathbb{R}^d$ is the joint configuration of the d -DOF robot hand, and $T \in SE(3)$ is the end-effector pose of the robot’s wrist consisting of translation and orientation. A *trial* is an instance of nonprehensile pushing or pulling, defined by a given direction $u_{\text{dir}} \in \mathbb{R}^3$ (with z-component of 0) resulting in the target object position as $u_{\text{targ}} \in \mathbb{R}^3$. The objective is to generate a hand pose H such that, if a motion planner moves the hand to H and then translates it along u_{dir} , the object moves closer to the target u_{targ} . The object’s distance to u_{targ} must be below a threshold for a trial to be considered a success.

4 Method

DexNoMa consists of the following steps. First, we generate a large dataset of hand poses for nonprehensile pushing and pulling (Sec. 4.1). Second, we use this data to train a diffusion model to synthesize hand poses conditioned on object geometry (Sec. 4.2). Third, during deployment, we generate hand poses and perform motion planning to do the pushing or pulling (Sec. 4.3).

4.1 Dataset Generation for Dexterous Nonprehensile Pushing and Pulling

We first generate diverse hand poses for pushing and pulling various objects in simulation. To do this, we take inspiration from prior work on generating diverse hand poses for *grasping* [15, 16, 17, 18, 20, 50] by casting the hand synthesis problem as minimizing an energy function via optimization [37]. Unlike those works, our focus is on pushing and pulling actions instead of

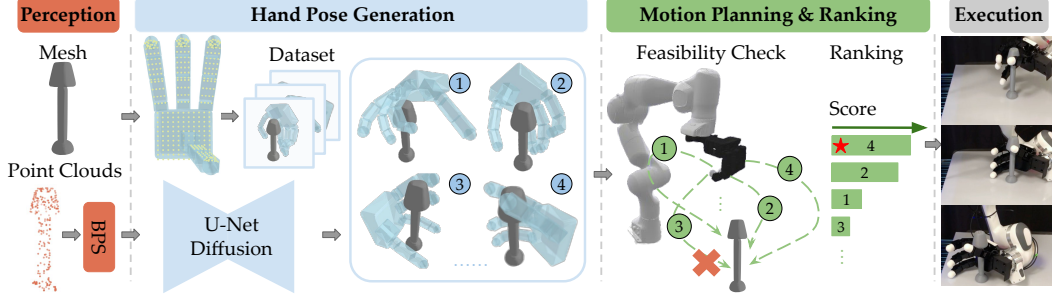


Figure 2: Overview of DexNoMa. We present a large-scale dataset of hand poses specifically for pushing or pulling, and leverage it to train a diffusion model. During execution time, given an object, we obtain its basis point set representation [22] and pass that to our trained diffusion model, which uses the architecture from [38]. This model synthesizes diverse floating pre-contact hand poses formed from our large-scale data generation pipeline (Sec. 4.1). Given these hand poses, we then check their feasibility in a physics simulator by adding the arm back in and performing motion planning [57]. We rank the feasible hand poses (e.g., “3” is infeasible in the example here) and select the best performing one (e.g., “4” in our example) and execute it in the real world.

119 grasping. To enable optimization, we first define a set of candidate contact points sampled across
 120 the hand surface. Different regions of the hand have different candidate points to encourage broad
 121 contact across the palm and fingers. For the palm and finger (excluding fingertips) regions, we sam-
 122 ple points uniformly over the rigid body surface. For the fingertips, we sample from a denser set of
 123 points uniformly on the unit hemisphere for each tip. See the Appendix for details of the distribution
 124 of candidate contact points (Figure 10 and Table 2).

125 With the sampled contact point candidates, we run an optimization algorithm following the sampling
 126 strategy from [15, 16] that iteratively minimizes an energy function E to generate hand poses. We
 127 adapt the energy function from [15] to better suit our nonprehensile manipulation tasks, resulting in:

$$E = E_{fc} + w_{dis}E_{dis} + w_{joints}E_{joints} + w_{pen}E_{pen} + w_{dir}E_{dir} + w_{arm}E_{arm}, \quad (1)$$

128 where E_{fc} is a force closure estimator [37], E_{dis} penalizes hand-to-object distance (thus encouraging
 129 proximity), E_{joints} penalizes joint violations, and E_{pen} penalizes penetration between hand-object,
 130 hand-table and hand self-collision contacts. See [15, 16] for further details. The w terms are all
 131 scalar coefficients; we adopt the values from prior work and tune the weights (see the Appendix) for
 132 the following two new terms. To adapt the energy from Eq. 1 to pushing or pulling in a particular
 133 direction $u_{dir} \in \mathbb{R}^3$, we introduce E_{dir} and E_{arm} , which use the normal vector of the palm $v_{palm} \in$
 134 \mathbb{R}^3 . The E_{dir} term encourages v_{palm} to align with u_{dir} , and E_{arm} encourages hand poses that are
 135 kinematically feasible when attached to the robot arm. Formally, we define E_{dir} and E_{arm} as:

$$E_{dir} = -\frac{u_{dir}^T v_{palm}}{\|u_{dir}\|_2 \|v_{palm}\|_2} \quad \text{and} \quad E_{arm} = \max(0, (v_{palm})_z) \quad (2)$$

136 where $(v_{palm})_z$ is the z -component of the palm’s normal vector (in the world frame). Intuitively,
 137 aligning u_{dir} and v_{palm} promotes more stable object-palm directional contact. Furthermore, if the
 138 palm faces upwards, then the rest of the arm must be below it. Thus, it is likely to lead to an
 139 infeasible robot configuration due to robot-table intersections, so E_{arm} is nonzero (i.e., worse). To
 140 inject randomness (and thus diversity) in the sampling process, we randomly resample a subset of
 141 the contact point indices from the set of valid candidates (Figure 10) when generating a new hand
 142 pose. We use RMSProp [58] to update translation, rotation and joint angles with step size decay,
 143 then minimize the energy function with Simulated Annealing [59] to adjust parameters.

144 **Hand Pose Validation in Simulation.** After optimizing contact points to generate candidate hand
 145 poses, we must *validate* whether they can lead to successful pushing or pulling. To do this, we use
 146 IsaacGym [21], a GPU-accelerated physics simulator that has been used in prior work for filtering
 147 grasp poses [15, 16]. We define a push or a pull as successful if, after executing a 20 cm translation,
 148 the object’s center is within 3 cm of the target position *and* the object’s orientation changes by no
 149 more than 45 degrees relative to its original configuration. The optimization process has a low

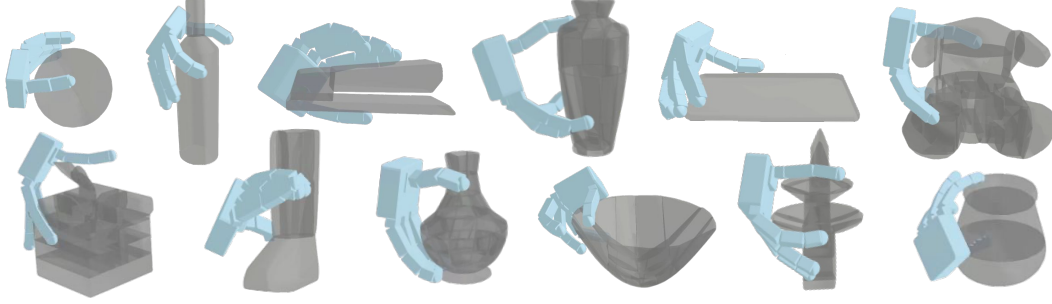


Figure 3: Examples of nonprehensile hand pushing poses from optimizing our energy function (Eq. 1). These have all been validated in IsaacGym simulation. In all examples, the intended object pushing direction is to the right. These data points are used to train our diffusion model (see Sec. 4.2).

150 success rate because it does not account for the full dynamics of pushing and pulling. Thus, we
 151 augment successful hand poses by adding slight noise to the pose parameters. We get 10X more
 152 augmented hand poses (from initially successful poses), and over the whole augmented dataset,
 153 58% are successful. From extensive parallel experiments, we generate a dataset containing 2,391
 154 objects with 1,387,632 successful hand poses. See the Appendix for more details.

155 4.2 Training a Diffusion Model to Predict Hand Poses

156 To generate hand poses, we adapt a conditional U-Net [60] from the diffusion policy architec-
 157 ture [11], and train it with the Denoising Diffusion Probabilistic Models (DDPM) objective [61].
 158 Diffusion models are well-suited for this task as they can learn complex, high-dimensional distribu-
 159 tions. The forward process gradually adds Gaussian noise to the hand configuration H , while the
 160 reverse process reconstructs the original pose H by iteratively denoising conditioned on the object’s
 161 geometry. The model is trained to minimize denoising error. To represent the observation, we use a
 162 4096-dimensional Basis Point Set (BPS) [22] representation $B \in \mathbb{R}^{4096}$ based on the object’s point
 163 cloud P . This representation, which is also used in [15, 38], encodes each object as a fixed-length
 164 vector of shortest distances between canonical basis points and the points in P . BPS captures geo-
 165 metric properties in a compact manner and simplifies the design of the diffusion model. Given this
 166 trained diffusion model, at test time it can be used to generate diverse hand poses which we can
 167 select for motion planning. See Figure 2 and Appendix A.2 for more information.

168 4.3 Arm-Hand Motion Planning and Evaluation

169 During deployment, the diffusion model generates candidate hand poses. We then integrate the
 170 Franka arm into full arm-hand motion planning to select hand poses which are kinematically fea-
 171 sible and avoid environment collisions, such as arm-table intersections (which are not considered
 172 in Sec. 4.1). See Figure 2 (right half) for an overview. Each hand pose $H = (\theta, T)$ is initially
 173 expressed in the object frame. We use the object’s initial configuration S_{obj} and intended direction
 174 u_{dir} to transform H to the world frame, and supply that to the cuRobo planner [57] to generate a
 175 complete motion plan for the Franka arm. In this process, we discard infeasible trajectories (and
 176 thus, the associated hand poses) to only keep the feasible arm-hand trajectories. To select which of
 177 the feasible trajectories to execute, we associate each with a custom analytical score V , defined as:

$$V(H = (\theta, T)) = \alpha L_{\text{goal}} + \beta L_{\text{coll}} + \gamma L_{\text{dir}}, \quad (3)$$

178 where L_{goal} measures the Euclidean distance between the object’s final position and the target po-
 179 sition, L_{coll} indicates whether a collision occurred during execution (1 if a collision occurs, 0 other-
 180 wise), and L_{dir} encourages the palm’s orientation to align with the pushing direction. For L_{dir} , we
 181 set it equal to the E_{dir} term from the energy function (Eq. 1). The α , β , and γ are hyperparameters.

182 **Multi-step Planning.** While we mainly study DexNoMa for single open-loop pushes (or pulls) to
 183 targets, our framework naturally extends to multi-step planning. In scenarios with obstacles, we
 184 first compute a collision-free global path using RRT* [62]. Then, we sequentially plan hand poses



Figure 4: The objects we use in our real-world nonprehensile manipulation experiments, including 3D printed and common (“Daily”) objects. See Sec. 5.2 for more details.

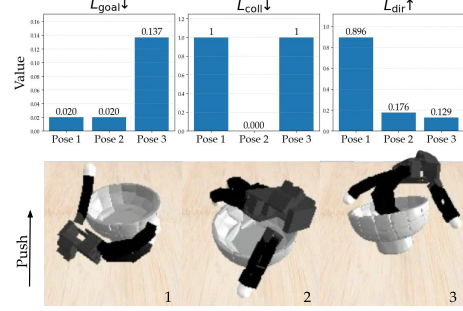


Figure 5: Visualization of L_{goal} , L_{coll} , and L_{dir} values in $V(H)$ from Eq. 3 on three simulated hand poses. See Sec. 5.3 for more details.

185 to reach each intermediate waypoint. Given an object, the same hand pose may be feasible only
 186 in certain pushing or pulling directions due to robot and hand kinematics. The waypoints from
 187 RRT* may require planning pushes across challenging directions, which highlights the importance
 188 of generating diverse hand poses for varying object positions and directions.

189 5 Experiments

190 Through simulation and real-world experiments, we aim to investigate the following questions: (1)
 191 Can we learn feasible and effective hand poses from our large-scale dataset? (2) Is DexNoMa robust
 192 to different pushing and pulling directions for visually diverse objects? (3) Can DexNoMa serve as
 193 a reliable module for downstream manipulation tasks such as multi-step pushing around obstacles?

194 5.1 Simulation Experiments and Results

195 We evaluate the quality of the hand pose generation pipeline using
 196 IsaacGym [21]. To quantify the effectiveness of our trained model
 197 and dataset, we report the number of successfully pushed objects as a
 198 function of training data size. We train our diffusion model on vary-
 199 ing subsets of the full dataset (of 1.3M hand poses) and evaluate on
 200 300 unseen objects from the test set. For each test object, we sam-
 201 ple 200 candidate hand poses. An object is considered “successful”
 202 if at least one feasible hand pose results in success. Table 1 reports
 203 results over 3 different seeds, which shows that our model generates
 204 feasible pushing poses more reliably with larger training sets, which validates large-scale supervi-
 205 sion. The growth is not strictly linear, suggesting room for improvement via better model tuning
 206 or data strategies. Qualitatively, our generated hand poses are diverse across object geometries and
 207 exhibit pushing intent (see the Appendix for more discussion). A common failure mode is that some
 208 poses still collide with the object, which motivates the inclusion of the collision term in Eq. 3.

Data Size	# of Objects
2%	41.67 \pm 10.21
20%	102.67 \pm 5.85
50%	110.33 \pm 29.67
100%	169.33 \pm 15.18

Table 1: Number of objects with at least one feasible pushing hand pose out of 300.

209 5.2 Real-World Experiments

210 We evaluate DexNoMa on a real robot to check if our nonprehensile hand poses successfully trans-
 211 fer to reality. Our hardware setup consists of a Franka Panda arm equipped with a four-finger,
 212 16-DOF Allegro Hand (see the Appendix). It operates over a tabletop cutting board with dimen-
 213 sions 60cm \times 60cm. We use a mix of objects, including 3D-printed and common (“Daily”) items
 214 (shown in Figure 4). All evaluation objects are unseen during training. For 3D-printed objects, we
 215 use their known meshes to directly compute their BPS representation. For the other objects, we
 216 follow the pipeline proposed in [63] to obtain real-world object point clouds (and thus, the BPS).
 217 We reconstruct object meshes by using Nerfstudio [64] to compute COLMAP reconstructions [65].
 218 We also use Stable Normal [66] to generate normal maps. Then, we employ 2D Gaussian Splat-
 219 ting [67] to obtain the point clouds. While this reconstruction pipeline introduces some noise, it is

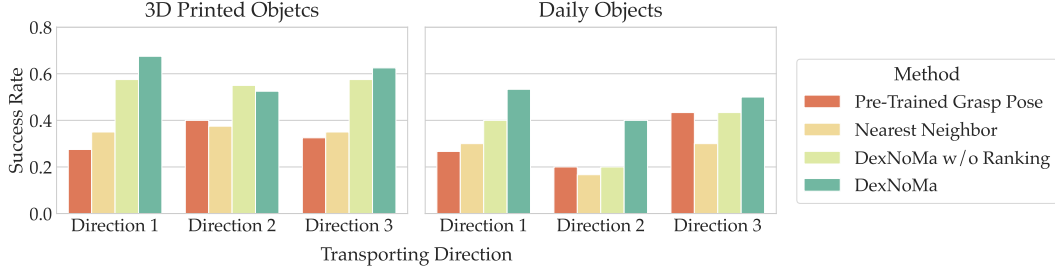


Figure 6: Nonprehensile manipulation success rates from DexNoMa and baselines, across different 3D printed (left) and daily objects (right), and with three directions evaluated. Each bar aggregates success rates from 40 trials (left bar plot) and 30 trials (right bar plot). See Sec. 5.2 and 5.3 for more details.

sufficient for DexNoMa to predict effective hand poses. In contrast, we empirically observed that optimization-based methods are more sensitive to mesh quality and often fail under these conditions.

Baselines and Ablations. We compare DexNoMa with the following methods.

- **Pre-Trained Grasp Pose:** We use a pre-trained grasp synthesis model from Lum et al. [15] using NeRF [68]. For each object, we train a NeRF representation, then query their pre-trained model for a grasp. This evaluates how well a grasping-centric model generalizes to nonprehensile tasks.
- **Nearest Neighbor (NN):** Given a test object, we find the training object with the most similar BPS representation (in terms of Euclidean distance) and retrieve its associated hand poses. We then do the same motion planning pipeline as in DexNoMa. This tests out-of-distribution generalization with a retrieval-only approach compared to our proposed generative model.
- **DexNoMa w/o Ranking:** An ablation that excludes analytical ranking of hand poses (ignores Eq. 3) and executes a random feasible pose. This tests the usefulness of Eq. 3 in selecting poses.

Experiment Protocol and Evaluation. For each object, we test three pushing directions uniformly distributed around a circle. Along each direction, the robot executes the hand pose and planned motion five times, all with a fixed push length of 20 cm. A human manually places the object in a relatively consistent pose between trials. A trial is successful if the object’s center is within 3 cm of the target position, the hand maintains contact throughout, and it does not lead to task failure modes such as toppling or loss of control. For NN and DexNoMa w/o Ranking, we randomly sample hand poses among the feasible planned actions. For Pre-trained Grasp Pose, we execute the best actions from its output. For our method, we execute the one with the highest analytical score from Eq. 3.

5.3 Real-World Results

We summarize quantitative results in Figure 6, which shows that DexNoMa outperforms or matches alternative methods for both object categories. As shown in Figure 7, the **Pre-Trained Grasp Pose** baseline suffers from two major issues. First, the hand pose is not conditioned on the pushing direction, which means during the push, the object is likely to slide off the hand due to limited support (Figure 7, second row). Second, some objects are unsuitable for grasping due to their geometry or awkward aspect ratios (e.g., a flat box with limited area for enclosure). Additionally, the similarity-based **Nearest Neighbor** baseline struggles due to limited granularity in object geometry matching, motivating the need for our geometry-conditioned generative model. For **DexNoMa w/o ranking**, we observe that its hand poses are more likely to collide with the table or objects. To further investigate this ablation, Figure 5 shows three different hand poses. The first one has a low collision score because it is easy to collide with the table, while the third collides with the objects and scores low on the palm direction. The second hand pose leads to a successful push in real-world experiments. This suggests the importance of our ranking system via Eq. 3. **DexNoMa** outperforms baselines in all directions tested in Figure 6, demonstrating the robustness of its generated hand poses for nonprehensile manipulation. Figure 7 (first row) demonstrates using the palm and thumb to provide strong support moving the object forward, and the third row shows using the thumb and index finger to form a circular shape support for the thinner upper parts of the object while providing force at the bottom, aiding stable movement. For more rollouts, see the Appendix and the website.

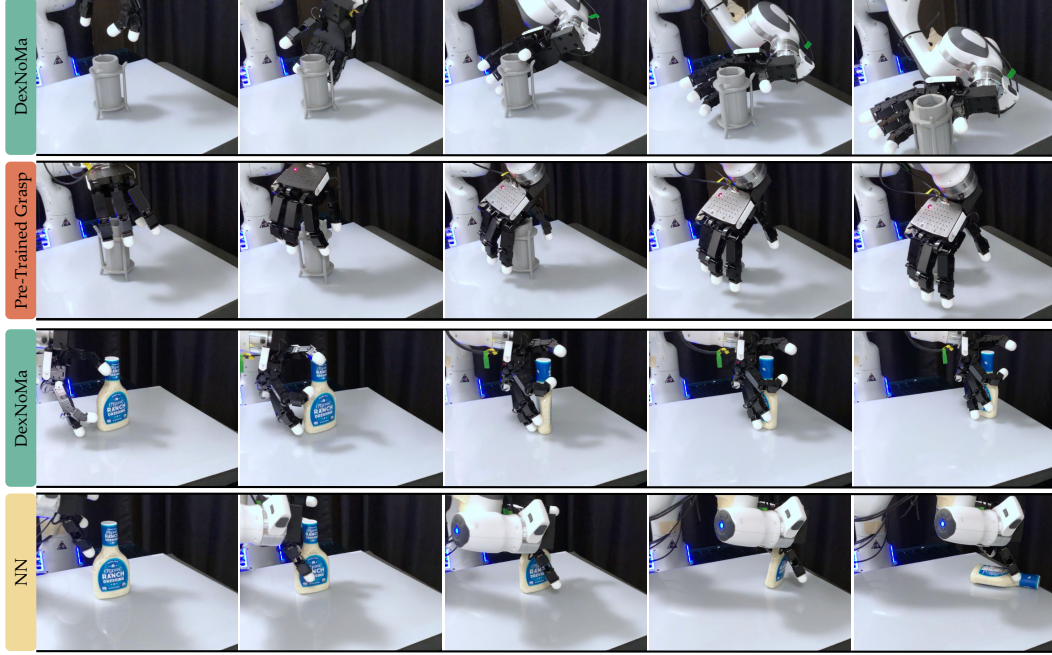


Figure 7: Comparison between DexNoMa and baselines. The first two rows show DexNoMa (success) and Pre-Trained Grasp (failure) while pushing a 3D-printed vase forward (i.e., away from the robot). The last two rows show DexNoMa (success) and NN (failure) while pushing a ranch bottle to the right.

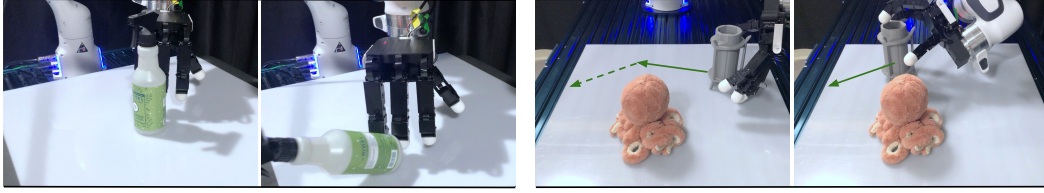


Figure 8: Example of a typical failure case using the Fixed Hand Pose strategy, which topples the spray.

Figure 9: Example of multi-step pushes using DexNoMa, which avoids the central obstacle.

259 *Fixed Hand Pose:* Inspired by prior pushing work [12], we manually define a “spatula” hand pose
 260 with the fingers spread flat (see Figure 8) to assess whether simple flat-hand strategies suffice for
 261 diverse objects. We perform a case study on the 6 objects in Figure 4 that are taller than 20 cm.
 262 We push each object 10 times, with 5 pushes for each of 2 directions, (the third direction results in
 263 kinematic errors). We get a relatively low 18/60 success rate, suggesting insufficient object support.

264 **Multi-step Planning.** Selecting a kinematically feasible hand pose for a given object state S_{obj}
 265 and direction u_{dir} is challenging in multi-step planning, as different waypoints may require different
 266 hand poses. Our method resolves this by identifying valid poses across object configurations and
 267 coupling pose selection with kinematic feasibility (see Sec. 4.3). By doing so, DexNoMa can be
 268 used to perform multiple pushes. Figure 9 shows a multi-step pushing sequence using DexNoMa.
 269 The robot uses two different hand poses to push the 3D-printed vase, as the first hand pose may not
 270 be ideal for the second hand pose, which shows the benefit of re-planning.

271 6 Conclusion

272 In this work, we propose DexNoMa, a dataset and method for nonprehensile object pushing and
 273 pulling using a high-DOF Allegro Hand. Our extensive real-world results show that DexNoMa en-
 274 ables diverse and effective pre-contact hand poses for different combinations of objects and pushing
 275 directions. We also demonstrate its usage for multi-step planning. We hope that this inspires future
 276 work on dexterous nonprehensile robotic manipulation.

7 Limitations

While promising, the DexNoMa approach has some limitations that motivate exciting directions for future work. First, it is difficult to get high success rates during the hand poses synthesis phase, and thus our method has room to improve for more data-efficient sampling. Second, we do not consider orientation when we evaluate pushing or pulling in the real world, drawing an incomplete picture of performance. Third, we only study pushing and pulling as examples of nonprehensile manipulation, which does not exhaustively characterize all possible nonprehensile manipulation procedures. Finally, it would be an interesting next step to make nonprehensile pushing or pulling truly closed-loop so it can react in real-time to unexpected disturbances such as object toppling.

References

- [1] M. T. Mason. Mechanics and Planning of Manipulator Pushing Operations. In *International Journal of Robotics Research (IJRR)*, 1986.
- [2] K. M. Lynch. *Nonprehensile Robotic Manipulation: Controllability and Planning*. PhD thesis, Carnegie Mellon University, The Robotics Institute, 1996.
- [3] M. T. Mason. Progress in Nonprehensile Manipulation. In *International Journal of Robotics Research (IJRR)*, 1999.
- [4] K. M. Lynch and M. T. Mason. Dynamic nonprehensile manipulation: Controllability, planning, and experiments. In *International Journal of Robotics Research (IJRR)*, 1999.
- [5] W. Zhou and D. Held. Learning to Grasp the Ungraspable with Emergent Extrinsic Dexterity. In *Conference on Robot Learning (CoRL)*, 2022.
- [6] W. Zhou, B. Jiang, F. Yang, C. Paxton, and D. Held. HACMan: Learning Hybrid Actor-Critic Maps for 6D Non-Prehensile Manipulation. In *Conference on Robot Learning (CoRL)*, 2023.
- [7] B. Jiang, Y. Wu, W. Zhou, C. Paxton, and D. Held. Hacman++: Spatially-grounded motion primitives for manipulation. In *Robotics: Science and Systems (RSS)*, 2024.
- [8] Y. Cho, J. Han, Y. Cho, and B. Kim. Corn: Contact-based object representation for non-prehensile manipulation of general unseen objects. In *International Conference on Learning Representations (ICLR)*, 2024.
- [9] J. Lyu, Z. Li, X. Shi, C. Xu, Y. Wang, and H. Wang. Dywa: Dynamics-adaptive world action model for generalizable non-prehensile manipulation. *arXiv preprint arXiv:2503.16806*, 2025.
- [10] K. Zhang, B. Li, K. Hauser, and Y. Li. Adaptigraph: Material-adaptive graph-based neural dynamics for robotic manipulation. In *Robotics: Science and Systems (RSS)*, 2024.
- [11] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion Policy: Visuomotor Policy Learning via Action Diffusion. In *Robotics: Science and Systems (RSS)*, 2023.
- [12] Y. Wang, Y. Li, K. Driggs-Campbell, L. Fei-Fei, and J. Wu. Dynamic-resolution model learning for object pile manipulation. In *Robotics: Science and Systems (RSS)*, 2023.
- [13] K. Shaw, A. Agarwal, and D. Pathak. LEAP Hand: Low-Cost, Efficient, and Anthropomorphic Hand for Robot Learning. In *Robotics: Science and Systems (RSS)*, 2023.
- [14] Y. Wang, Y. Li, Y. Yang, and Y. Chen. Dexterous non-prehensile manipulation for ungraspable object via extrinsic dexterity. *arXiv preprint arXiv:2503.23120*, 2025.
- [15] T. G. W. Lum, A. H. Li, P. Culbertson, K. Srinivasan, A. D. Ames, M. Schwager, and J. Bohg. Get a Grip: Multi-Finger Grasp Evaluation at Scale Enables Robust Sim-to-Real Transfer. In *Conference on Robot Learning (CoRL)*, 2024.

- [16] R. Wang, J. Zhang, J. Chen, Y. Xu, P. Li, T. Liu, and H. Wang. DexGraspNet: A large-scale robotic dexterous grasp dataset for general objects based on simulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [17] J. Zhang, H. Liu, D. Li, X. Yu, H. Geng, Y. Ding, J. Chen, and H. Wang. DexGraspNet 2.0: Learning Generative Dexterous Grasping in Large-scale Synthetic Cluttered Scenes. In *Conference on Robot Learning (CoRL)*, 2024.
- [18] Y. Xu, W. Wan, J. Zhang, H. Liu, Z. Shan, H. Shen, R. Wang, H. Geng, Y. Weng, J. Chen, et al. UniDexGrasp: Universal robotic dexterous grasping via learning diverse proposal generation and goal-conditioned policy. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [19] W. Wan, H. Geng, Y. Liu, Z. Shan, Y. Yang, L. Yi, and H. Wang. UniDexGrasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- [20] Y. Li, B. Liu, Y. Geng, P. Li, Y. Yang, Y. Zhu, T. Liu, and S. Huang. Grasp multiple objects with one hand. In *IEEE Robotics and Automation Letters (RA-L)*, 2024.
- [21] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, and G. State. Isaac Gym: High Performance GPU-Based Physics Simulation For Robot Learning. *arXiv preprint arXiv:2108.10470*, 2021.
- [22] S. Prokudin, C. Lassner, and J. Romero. Efficient learning on point clouds with basis point sets. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [23] J. Moura, T. Stouraitis, and S. Vijayakumar. Non-prehensile planar manipulation via trajectory optimization with complementarity constraints. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2022.
- [24] N. Chavan-Dafle, A. Rodriguez, R. Paolini, B. Tang, S. Srinivasa, M. Erdmann, M. T. Mason, I. Lundberg, H. Staab, and T. Fuhlbrigge. Extrinsic dexterity: In-hand manipulation with external forces. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2014.
- [25] W. Yang and M. Posa. Dynamic on-palm manipulation via controlled sliding. In *Robotics: Science and Systems (RSS)*, 2024.
- [26] S. Chen, A. Wu, and C. K. Liu. Synthesizing dexterous nonprehensile pregrasp for ungraspable objects. In *ACM SIGGRAPH*, 2023.
- [27] A. Wu, R. Wang, S. Chen, C. Eppner, and C. K. Liu. One-Shot Transfer of Long-Horizon Extrinsic Manipulation Through Contact Retargeting. *arXiv preprint arXiv:2404.07468*, 2024.
- [28] W. C. Agboh, J. Ichnowski, K. Goldberg, and M. R. Dogar. Multi-object grasping in the plane. In *International Symposium on Robotics Research (ISRR)*, 2022.
- [29] T. Yonemaru, W. Wan, T. Nishimura, and K. Harada. Learning to Group and Grasp Multiple Objects. *arXiv preprint arXiv:2502.08452*, 2025.
- [30] J. J. Liu, Y. Li, K. Shaw, T. Tao, R. Salakhutdinov, and D. Pathak. Factr: Force-attending curriculum training for contact-rich policy learning. In *Robotics: Science and Systems (RSS)*, 2025.
- [31] M. Liu, Z. Pan, K. Xu, K. Ganguly, and D. Manocha. Deep differentiable grasp planner for high-dof grippers. In *Robotics: Science and Systems (RSS)*, 2020.
- [32] A. T. Miller and P. K. Allen. Graspit! a versatile simulator for robotic grasping. *IEEE Robotics & Automation Magazine*, 2004.

- [33] H. Jiang, S. Liu, J. Wang, and X. Wang. Hand-object contact consistency reasoning for human grasps generation. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [34] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations (ICLR)*, 2014.
- [35] D. Turpin, L. Wang, E. Heiden, Y.-C. Chen, M. Macklin, S. Tsogkas, S. Dickinson, and A. Garg. Grasp’d: Differentiable contact-rich grasp synthesis for multi-fingered hands. In *European Conference on Computer Vision (ECCV)*, 2022.
- [36] D. Turpin, T. Zhong, S. Zhang, G. Zhu, E. Heiden, M. Macklin, S. Tsogkas, S. Dickinson, and A. Garg. Fast-grasp’d: Dexterous multi-finger grasp generation through differentiable simulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [37] T. Liu, Z. Liu, Z. Jiao, Y. Zhu, and S.-C. Zhu. Synthesizing diverse and physically stable grasps with arbitrary hand structures using differentiable force closure estimator. In *IEEE Robotics and Automation Letters (RA-L)*, 2022.
- [38] Z. Weng, H. Lu, D. Kragic, and J. Lundell. Dexdiffuser: Generating dexterous grasps with diffusion models. In *IEEE Robotics and Automation Letters (RA-L)*, 2024.
- [39] J. Bohg, A. Morales, T. Asfour, and D. Kragic. Data-driven grasp synthesis—a survey. In *IEEE Transactions on Robotics (T-RO)*, 2014.
- [40] O. Kroemer, S. Niekum, and G. Konidaris. A Review of Robot Learning for Manipulation: Challenges, Representations, and Algorithms. In *Journal of Machine Learning Research (JMLR)*, 2021.
- [41] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning Fine-Grained Bimanual Manipulation with Low-Cost Hardware. In *Robotics: Science and Systems (RSS)*, 2023.
- [42] Z. Fu, T. Z. Zhao, and C. Finn. Mobile ALOHA: Learning Bimanual Mobile Manipulation with Low-Cost Whole-Body Teleoperation. In *Conference on Robot Learning (CoRL)*, 2024.
- [43] H. Qi, A. Kumar, R. Calandra, Y. Ma, and J. Malik. In-Hand Object Rotation via Rapid Motor Adaptation. In *Conference on Robot Learning (CoRL)*, 2022.
- [44] J. Wang, Y. Yuan, H. Che, H. Qi, Y. Ma, J. Malik, and X. Wang. Lessons from Learning to Spin “Pens”. In *Conference on Robot Learning (CoRL)*, 2024.
- [45] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. Zhang. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- [46] OpenAI, M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba. Learning Dexterous In-Hand Manipulation. In *International Journal of Robotics Research (IJRR)*, 2019.
- [47] Y. Yuan, H. Che, Y. Qin, B. Huang, Z.-H. Yin, K.-W. Lee, Y. Wu, S.-C. Lim, and X. Wang. Robot Synesthesia: In-Hand Manipulation with Visuotactile Sensing. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [48] H. Jiang, Y. Wang, H. Zhou, and D. Seita. Learning to Singulate Objects in Packed Environments using a Dexterous Hand. In *International Symposium on Robotics Research (ISRR)*, 2024.

- [49] L. Xu, Z. Liu, Z. Gui, J. Guo, Z. Jiang, Z. Xu, C. Gao, and L. Shao. Dexsingrasp: Learning a unified policy for dexterous object singulation and grasping in cluttered environments. *arXiv preprint arXiv:2504.04516*, 2025.
- [50] S. He, Z. Shangguan, K. Wang, Y. Gu, Y. Fu, Y. Fu, and D. Seita. Sequential multi-object grasping with one dexterous hand. *arXiv preprint arXiv:2503.09078*, 2025.
- [51] K. Yao and A. Billard. Exploiting kinematic redundancy for robotic grasping of multiple objects. In *IEEE Transactions on Robotics (T-RO)*, 2023.
- [52] T. Chen, E. Cousineau, N. Kuppawamy, and P. Agrawal. Vegetable Peeling: A Case Study in Constrained Dexterous Manipulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2025.
- [53] T. Lin, Z.-H. Yin, H. Qi, P. Abbeel, and J. Malik. Twisting Lids Off with Two Hands. In *Conference on Robot Learning (CoRL)*, 2024.
- [54] Y. Ze, G. Zhang, K. Zhang, C. Hu, M. Wang, and H. Xu. 3d diffusion policy: Generalizable visuomotor policy learning via simple 3d representations. In *Robotics: Science and Systems (RSS)*, 2024.
- [55] T. G. W. Lum, O. Y. Lee, C. K. Liu, and J. Bohg. Crossing the Human-Robot Embodiment Gap with Sim-to-Real RL using One Human Demonstration. *arXiv preprint arXiv:2504.12609*, 2025.
- [56] J. Chen, Y. Chen, J. Zhang, and H. Wang. Task-oriented dexterous hand pose synthesis using differentiable grasp wrench boundary estimator. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2024.
- [57] B. Sundaralingam, S. K. S. Hari, A. Fishman, C. Garrett, K. V. Wyk, V. Blukis, A. Millane, H. Oleynikova, A. Handa, F. Ramos, N. Ratliff, and D. Fox. curobo: Parallelized collision-free minimum-jerk robot motion generation. *arXiv preprint arXiv:2310.17274*, 2023.
- [58] T. Tieleman and G. Hinton. Lecture 6.5—rmsprop: Divide the gradient by a running average of its recent magnitude, 2012. COURSERA: Neural Networks for Machine Learning.
- [59] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, 1983.
- [60] O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015.
- [61] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *Neural Information Processing Systems (NeurIPS)*, 2020.
- [62] S. Karaman and E. Frazzoli. Sampling-based algorithms for optimal motion planning. *International Journal of Robotics Research (IJRR)*, 2011.
- [63] H. Lou, Y. Liu, Y. Pan, Y. Geng, J. Chen, W. Ma, C. Li, L. Wang, H. Feng, L. Shi, et al. Robogs: A physics consistent spatial-temporal model for robotic arm with hybrid representation. *arXiv preprint arXiv:2408.14873*, 2024.
- [64] M. Tancik, E. Weber, E. Ng, R. Li, B. Yi, J. Kerr, T. Wang, A. Kristoffersen, J. Austin, K. Salahi, A. Ahuja, D. McAllister, and A. Kanazawa. Nerfstudio: A modular framework for neural radiance field development. In *ACM SIGGRAPH 2023 Conference Proceedings, SIGGRAPH ’23*, 2023.
- [65] J. L. Schönberger and J.-M. Frahm. Structure-from-Motion Revisited. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

- 449 [66] C. Ye, L. Qiu, X. Gu, Q. Zuo, Y. Wu, Z. Dong, L. Bo, Y. Xiu, and X. Han. Stablenormal:
450 Reducing diffusion variance for stable and sharp normal. *arXiv preprint arXiv:2406.16864*,
451 2024.
- 452 [67] B. Huang, Z. Yu, A. Chen, A. Geiger, and S. Gao. 2d gaussian splatting for geometrically
453 accurate radiance fields. In *Special Interest Group on Computer Graphics and Interactive*
454 *Techniques Conference Conference Papers '24*, SIGGRAPH '24. ACM, 2024.
- 455 [68] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. NeRF:
456 Representing Scenes as Neural Radiance Fields for View Synthesis. In *European Conference*
457 *on Computer Vision (ECCV)*, 2020.
- 458 [69] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, et al. Grounding
459 dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint*
460 *arXiv:2303.05499*, 2023.
- 461 [70] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland,
462 L. Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint*
463 *arXiv:2408.00714*, 2024.
- 464 [71] T. Ren, S. Liu, A. Zeng, J. Lin, K. Li, H. Cao, J. Chen, X. Huang, Y. Chen, F. Yan, Z. Zeng,
465 H. Zhang, F. Li, J. Yang, H. Li, Q. Jiang, and L. Zhang. Grounded sam: Assembling open-
466 world models for diverse visual tasks, 2024.
- 467 [72] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. White-
468 head, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick. Segment Anything. *arXiv preprint*
469 *arXiv:2304.02643*, 2023.
- 470 [73] Q. Jiang, F. Li, Z. Zeng, T. Ren, S. Liu, and L. Zhang. T-rex2: Towards generic object detection
471 via text-visual prompt synergy, 2024.