

D³Fields: Dynamic 3D Descriptor Fields for Zero-Shot Generalizable Robotic Manipulation

Anonymous Author(s)

Affiliation

Address

email

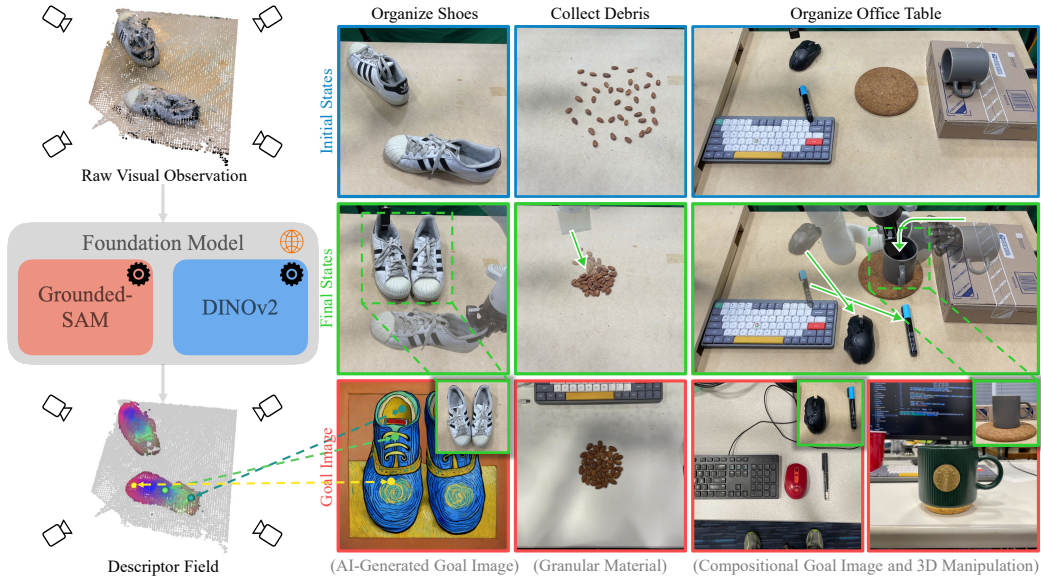


Figure 1: **D³Fields Representation and Application to Various Manipulation Tasks.** D³Fields take in multi-view RGBD images and encode semantic features and instance masks using foundational models. The gray and colored points in the bottom left visualize background and semantic features mapped to RGB space using Principal Component Analysis (PCA), demonstrating consistency across instances. We use our representation for diverse tasks in a zero-shot manner. These tasks are defined by 2D goal images with diverse instances and styles. We address pick-and-place tasks such as shoe organization and tasks requiring dynamic modeling like collecting debris. We also demonstrate in the office table organization that our framework can accomplish 3D manipulation and compositional task specification.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18

Abstract:

Scene representation has been a crucial design choice in robotic manipulation systems. An ideal representation should be 3D, dynamic, and semantic to meet the demands of diverse manipulation tasks. However, previous works often lack all three properties simultaneously. In this work, we introduce D³Fields — dynamic 3D descriptor fields. These fields capture the dynamics of the underlying 3D environment and encode both semantic features and instance masks. Specifically, we project arbitrary 3D points in the workspace onto multi-view 2D visual observations and interpolate features derived from foundational models. The resulting fused descriptor fields allow for flexible goal specifications using 2D images with varied contexts, styles, and instances. To evaluate the effectiveness of these descriptor fields, we apply our representation to a wide range of robotic manipulation tasks in a zero-shot manner. Through extensive evaluation in both real-world scenarios and simulations, we demonstrate that D³Fields are both generalizable and effective for zero-shot robotic manipulation tasks. In quantitative comparisons with state-of-the-art dense descriptors, such as Dense Object Nets and DINO, D³Fields exhibit significantly better generalization abilities and manipulation accuracy.

19 1 Introduction

20 The choice of scene representation is critical in robotic systems. An ideal representation should be
21 simultaneously 3D, dynamic, and semantic to meet the needs of various robotic manipulation tasks
22 in our daily lives. However, previous research into scene representations in robotics often does not
23 encompass all three properties. Some representations exist in 3D space [1, 2, 3, 4], yet they overlook
24 semantic information. Others focus on dynamic modeling [5, 6, 7, 8], but only consider 2D data.
25 Some other works are limited by only considering semantic information such as object instance and
26 category [9, 10, 11, 12, 13].

27 In this work, we aim to satisfy all three criteria by introducing D³Fields, unified descriptor fields
28 that are 3D, dynamic, and semantic. D³Fields take in arbitrary points in the 3D world coordinate
29 frame and output both geometric and semantic information related to these points. This includes
30 the instance mask, dense semantic features, and the signed distance to the object surface. Notably,
31 deriving these descriptor fields requires no training and is conducted in a zero-shot manner using
32 large foundational vision models and vision-language models (VLMs). Specifically, we first use
33 Grounding-DINO [14], Segment Anything (SAM) [15], XMem [16], and DINOv2 [17] to extract
34 information from multi-view 2D RGB images. We then project the 3D points back to each camera,
35 interpolate to compute representations from each view, and fuse these data to derive the descriptors
36 for the associated 3D points, as shown in Fig. 1 (left). By leveraging the dense semantic feature and
37 instance mask of our representation, we can robustly track 3D points of the target object instance
38 and train dynamics models. These learned dynamics models can then be incorporated into a Model-
39 Predictive Control (MPC) framework to plan for manipulation tasks.

40 Notably, the derived representations allow for goal specification using 2D images sourced from the
41 Internet, phones, or those generated by AI models. Such goal images have been challenging to
42 manage with previous methods, because they contain varied styles, contexts, and object instances
43 different from the robot’s workspace. Our proposed D³Fields can establish dense correspondences
44 between the robot workspace and the target configurations. These correspondences give us the task
45 objective, enabling us to plan the robot’s actions with the learned dynamics model within the MPC
46 framework. This task execution process does not require any further training, offering a flexible and
47 convenient interface for humans to instruct robots.

48 We evaluate our method across a wide range of household robotic manipulation tasks in a zero-
49 shot manner. These tasks include organizing shoes, collecting debris, and organizing office desks,
50 as shown in Fig. 1 (right). Furthermore, we offer detailed quantitative comparisons between our
51 method and other state-of-the-art dense descriptor techniques. Our results indicate that our approach
52 significantly outperforms in terms of generalizability and manipulation accuracy.

53 To summarize our contributions: (1) We introduce a novel representation, D³Fields, that is **3D**,
54 **dynamic**, and **semantic**. (2) We present a novel and flexible goal specification method using 2D
55 images that incorporate a range of styles, contexts, and instances. (3) Our proposed robotic manip-
56 ulation framework supports zero-shot generalizable manipulation applicable to a broad spectrum of
57 household tasks.

58 2 Related Works

59 2.1 Foundation Models for Robotics

60 Foundation models generally refer to those trained on broad data, often using self-supervision at
61 scale, which can then be adapted (e.g., fine-tuned) to various downstream tasks. Large Language
62 Models (LLMs) have showcased promising reasoning abilities for language. Robotics researchers
63 have recently released a series of works that leverage LLMs, including SayCan [18] and Inner Mono-
64 logue [19], to directly generate robot plans. Some later works have used LLMs as a code generator:
65 Code as Policies [20] uses 2D object detectors as the perception API, whereas VoxPoser [21] cre-
66 ates a 3D value map. Yet, their perception modules fall short in modeling the precise geometry and

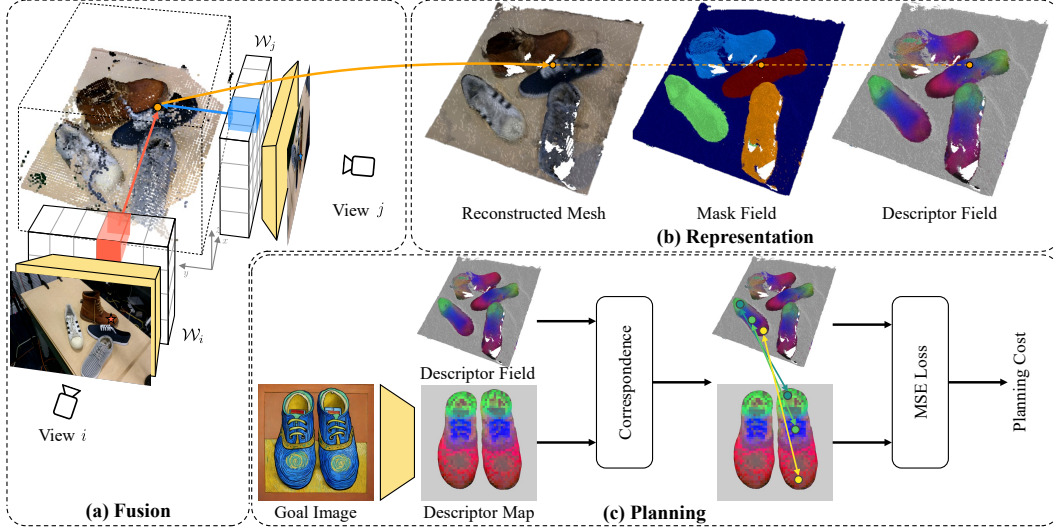


Figure 2: **Overview of the proposed framework.** (a) The fusion process fuses RGBD observations from multiple views. Each view is processed by foundation models to obtain the feature volume \mathcal{W} . Arbitrary 3D points are processed through projection and interpolation. (b) After fusing information from multiple views, we obtain an implicit distance function to reconstruct the mesh form. We also have instance masks and semantic features for evaluated 3D points, as shown by the mask field and descriptor field in the top right subfigure. (c) Given a 2D goal image, we use foundation models to extract the descriptor map. Then we correspond 3D features to 2D features and define the planning cost based on the correspondence.

67 dynamics of objects. Our D³Fields aim to address this by focusing on detailed 3D geometry and
 68 dynamics.

69 Meanwhile, foundational vision models, such as SAM [15] and DINOv2 [17], have demonstrated
 70 impressive zero-shot generalization capabilities across various vision tasks. However, their focus
 71 is primarily on 2D vision tasks. Grounding these models in a dynamic 3D environment remains a
 72 challenge. The recent GROOT project showcases how to construct 3D object-centric representa-
 73 tions using foundational models and exhibits notable few-shot generalization capabilities [22]. Still,
 74 GROOT does not emphasize learning about object dynamics or achieving zero-shot generalizable
 75 robotic manipulation.

76 2.2 Representation for Visual Robotic Manipulation

77 Scene representation has been a pivotal component in robotic manipulation systems. Some early
 78 work relies on 2D representations, such as bounding boxes [23, 24]. Many recent methods construct
 79 particle representations of the environment and employ learned dynamics to capture the system’s
 80 underlying structure [25, 3, 7, 8, 26, 27, 28, 29]. They demonstrate impressive results in unstructured
 81 environments and with non-rigid objects. However, they are not semantic, which can hinder their
 82 ability to generalize to new tasks and scenarios. Some research opts for a fixed-dimension latent
 83 vector derived from high-dimensional sensory inputs as the representation [30, 5, 6, 31, 32, 33, 34,
 84 35, 36, 2], but such a representation does not scale well to complex manipulation tasks that require
 85 high precision and explicit scene structures. Other approaches use 6 DoF object poses as their
 86 representation [9, 10, 37, 38], though focusing primarily on grasping tasks instead of more dynamic
 87 ones. In this work, we aim to address these issues by introducing D³Fields, a representation that
 88 models dynamic 3D environments at varying semantic levels.

89 2.3 Neural Fields for Robotic Manipulation

90 Researchers have presented a variety of works using neural fields as a representation for robotic
 91 manipulation [39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 41]. Among them, Neural
 92 Descriptor Fields are the most relevant to ours [42]. They build neural feature fields that generalize to

93 different instances with several demonstrations; but they focus on learning geometric, not semantic
 94 features, which hinders cross-category generalization.

95 Recently, a series of works distilled neural feature fields using foundation models such as CLIP
 96 and DINO for supervision [53, 54]. LeRF distills neural feature fields to handle open-vocabulary
 97 3D queries and develops task-oriented grasping based on it [55, 56]. Shen et al. [57] use a similar
 98 distilled feature field for the grasping task. Both methods require dense camera views to train the
 99 neural field. GNFactor addresses this by introducing a voxel encoder [58]. However, distilling
 100 foundation models to create neural feature fields has drawbacks: (1) They often require dense camera
 101 views for a quality field. (2) Distilled neural fields need retraining for new scenes, limiting their
 102 generalization and making them ineffective for dynamic scenes. In contrast, our D³Fields do not
 103 need extra training for new scenes and can work with sparse views and dynamic settings.

104 3 Method

105 In this section, we introduce the problem formulation in Section 3.1 and define camera transforma-
 106 tion and projection notations in Section 3.2. The construction of D³Fields is detailed in Section 3.3.
 107 Section 3.4 discusses tracking keypoints and learning dynamics, while Section 4.3 showcases how
 108 our representation enables zero-shot generalizable manipulation skills.

109 3.1 Problem Formulation

110 Given a 2D goal image \mathcal{I} , we denote the corresponding scene representation as \mathbf{s}_{goal} . Our goal is to
 111 find the action sequence $\{a^t\}$ to minimize the task objective:

$$\begin{aligned} \min_{\{a_t\}} \quad & c(\mathbf{s}^T, \mathbf{s}_{\text{goal}}), \\ \text{s.t.} \quad & \mathbf{s}^t = g(\mathbf{o}^t), \quad \mathbf{s}^{t+1} = f(\mathbf{s}^t, a^t), \end{aligned} \quad (1)$$

112 where $c(\cdot, \cdot)$ is the cost function measuring the distance between the terminal representation \mathbf{s}^T and
 113 the goal representation \mathbf{s}_{goal} . Representation extraction function $g(\cdot)$ takes in the current multi-view
 114 RGBD observations \mathbf{o}^t and outputs the current representation \mathbf{s}^t . $f(\cdot, \cdot)$ is the dynamics function that
 115 predicts the future representation \mathbf{s}^{t+1} , conditioned on the current representation \mathbf{s}^t and action a^t .
 116 The optimization aims to find the action sequence $\{a_t\}$ that minimizes the cost function $c(\mathbf{s}^T, \mathbf{s}_{\text{goal}})$.

117 3.2 Notation: Camera Transformation and Projection

118 We assume all cameras' intrinsic parameters \mathbf{K} and extrinsic parameters \mathbf{T} are known. The camera
 119 i extrinsic parameters are defined as follows.

$$\mathbf{T}_i = \begin{bmatrix} \mathbf{R}_i & \mathbf{t}_i \\ \mathbf{0}^T & 1 \end{bmatrix} \in \mathbb{SE}(3), \quad (2)$$

120 where Euclidean group $\mathbb{SE}(3) := \{\mathbf{R}, \mathbf{t} \mid \mathbf{R} \in \mathbb{SO}^3, \mathbf{t} \in \mathbb{R}^3\}$. For a 3D point \mathbf{x} in the world frame,
 121 we could obtain projected pixel \mathbf{u}_i and distance to camera \mathbf{r}_i as follows:

$$\mathbf{u}_i = \pi(\mathbf{K}_i(\mathbf{R}_i\mathbf{x} + \mathbf{t}_i)), \quad \mathbf{r}_i = [0, 0, 1]^T(\mathbf{R}_i\mathbf{x} + \mathbf{t}_i), \quad (3)$$

122 where π performs perspective projection, mapping a 3D vector $p = [x, y, z]^T$ to a 2D vector $q =$
 123 $[x/z, y/z]^T$.

124 3.3 D³Fields Representation

125 We fuse observation \mathbf{o}^t from multiple views to build the implicit 3D descriptor fields $\mathcal{F}^t(\cdot)$. For
 126 simplicity, we will represent \mathbf{o}^t as \mathbf{o} , and $\mathcal{F}^t(\cdot)$ as $\mathcal{F}(\cdot)$ in this subsection. The implicit 3D descriptor
 127 field $\mathcal{F}(\cdot)$ is defined as

$$(\mathbf{d}, \mathbf{f}, \mathbf{p}) = \mathcal{F}(\mathbf{x}), \quad (4)$$

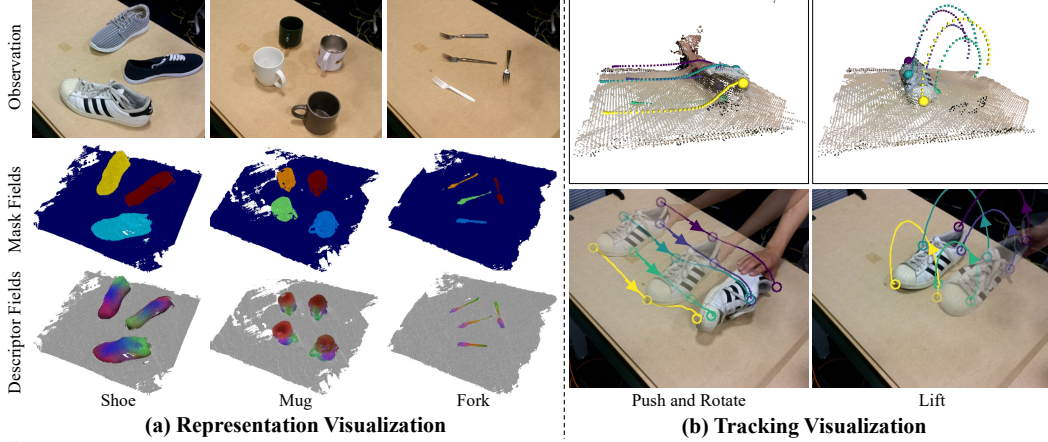


Figure 3: **Representation and Tracking Visualizations.** (a) To verify that the representation is both 3D and semantic, we visualize the representation across different object categories. Mask fields color 3D points based on their instance masks, which clearly differentiates between instances. Descriptor fields color 3D points by mapping features to RGB space using PCA. They display a consistent color pattern within a category, such as mug handles being colorized as green for different mug instances. (b) To demonstrate that our representation is dynamic, we apply it to tracking tasks and showcase two tracking examples, both of which involve 3D motions and partial observations in single views. The robust 3D tracking results confirm that our representation is 3D, dynamic, and semantic.

128 where \mathbf{x} is an arbitrary 3D point in the world frame, and $(\mathbf{d}, \mathbf{f}, \mathbf{p})$ is the corresponding geometric
 129 and semantic descriptor. $\mathbf{d} \in \mathbb{R}$ is the signed distance from \mathbf{x} to the surface. $\mathbf{f} \in \mathbb{R}^N$ represents the
 130 semantic information of N dimension. $\mathbf{p} \in \mathbb{R}^M$ denotes the instance probability distribution of M
 131 instances. M could be different across scenarios.

132 More specifically, we denote a single view RGBD observation from camera i as $\mathbf{o}_i = (\mathcal{I}_i, \mathcal{R}_i)$,
 133 where RGB image $\mathcal{I}_i \in \mathbb{R}^{H \times W \times 3}$, and depth image $\mathcal{R}_i \in \mathbb{R}^{H \times W}$. For an arbitrary 3D point \mathbf{x} , we
 134 project it to image space using Eq. 3 and use bilinear interpolation to obtain the corresponding depth
 135 $\mathbf{r}'_i = \mathcal{R}_i[\mathbf{u}_i]$. Then the descriptors from camera i are

$$\begin{aligned} \mathbf{d}_i &= \mathbf{r}'_i - \mathbf{r}_i, & \mathbf{d}'_i &= \max(\min(\mathbf{d}_i, \mu), -\mu), \\ \mathbf{f}_i &= \mathcal{W}_i^f[\mathbf{u}_i], & \mathbf{p}_i &= \mathcal{W}_i^p[\mathbf{u}_i], \end{aligned} \quad (5)$$

136 where DINOv2 [17] extracts the semantic feature volume $\mathcal{W}_i^f \in \mathbb{R}^{H \times W \times N}$ from RGB observa-
 137 tion \mathcal{I}_i . $\mathcal{W}_i^p \in \mathbb{R}^{H \times W \times M}$ is the instance mask volume using Grounded-SAM [14, 15]. μ is the
 138 truncation threshold for TSDF.

139 We fuse descriptors from all K views as follows:

$$v_i = H(\mathbf{d}_i + \mu), \quad w_i = \exp\left(\frac{\min(\mu - |\mathbf{d}_i|, 0)}{\mu}\right), \quad (6)$$

140 and then

$$\mathbf{d} = \frac{\sum_{i=1}^K v_i \mathbf{d}'_i}{\delta + \sum_{i=1}^K v_i}, \quad \mathbf{f} = \frac{\sum_{i=1}^K v_i w_i \mathbf{f}_i}{\delta + \sum_{i=1}^K v_i}, \quad \mathbf{m} = \frac{\sum_{i=1}^K v_i w_i \mathbf{m}_i}{\delta + \sum_{i=1}^K v_i}, \quad (7)$$

141 where H is the unit step function and δ is a small value to avoid numeric errors. $v_i = 0$ when
 142 \mathbf{x} is not observable in camera i , because if \mathbf{x} is occluded in camera i , it should not contribute to
 143 the descriptor of \mathbf{x} . In addition, we could only have a confident estimation when \mathbf{x} is close to the
 144 surface. Therefore, w_i will decay as $|\mathbf{d}_i|$ increases. For \mathbf{x} that is far away, \mathbf{f} and \mathbf{m} will degrade to
 145 0^T .

146 We convert the implicit field function $\mathcal{F}(\cdot)$ to a set of keypoints s . First, we create voxels $\mathbf{x} \in$
 147 $\mathbb{R}^{W \times L \times H \times 3}$ in the workspace and evaluate $(\mathbf{d}, \mathbf{f}, \mathbf{p}) = \mathcal{F}(\mathbf{x})$. We filter out $\mathbf{x}_i \in \mathbf{x}$ where \mathbf{d}_i is
 148 large or \mathbf{p}_i has a low probability to avoid empty space and the background. After obtaining filtered
 149 points \mathbf{x}' , we use farthest point sampling to find surface points $s \in \mathbb{R}^{3 \times n_s}$ of an instance.



Figure 4: **Qualitative results.** We qualitatively evaluate our proposed framework on household manipulation tasks, both in the real world and in simulation, encompassing tasks such as organizing utensils, fruits, shoes, food, and mugs. The figure highlights that our representation can generalize across varied instances, styles, and contexts. For instance, in the organizing fruits example, the goal image, unlike the workspace, is styled as a sketch drawing. Because our representation can map bananas with varied styles and appearances to similar features, the banana in the workspace can correspond to the banana in the sketch. This allows the task to be successfully completed. This wide range of tasks showcases the generalization capabilities and manipulation precision of our framework.

150 3.4 Keypoints Tracking and Dynamics Training

151 This section will present how to use the dynamic implicit 3D descriptor field $\mathcal{F}(\cdot)$ to track keypoints
 152 and train dynamics. Without losing generalization, consider the tracking of a single instance $s^t \in$
 153 $\mathbb{R}^{3 \times n_s}$. For clarity, we denote \mathbf{f} and \mathbf{d} from $\mathcal{F}(\cdot)$ as $\mathcal{F}_f(\cdot)$ and $\mathcal{F}_d(\cdot)$. We formulate the tracking
 154 problem as an optimization problem:

$$\min_{s^{t+1}} \|\mathcal{F}_f(s^{t+1}) - \mathcal{F}_f(s^0)\|_2. \quad (8)$$

155 Since $\mathcal{F}(\cdot)$ is differentiable, we could use a gradient-based optimizer. This method could be naturally
 156 extended to multiple-instance scenarios. We found that relying solely on features for tracking is
 157 unstable. We added rigid constraints and distance regularization for a more stable tracking.

158 Keypoint tracking enables dynamics model training on real data. We instantiate the dynamics model
 159 $f(\cdot, \cdot)$ as graph neural networks (GNNs). We follow [59] to predict object dynamics. Please refer
 160 to [25, 59] for more details on how to train the GNN-based dynamics model. The trained dynamics
 161 will be used for trajectory optimization in Section 3.5.

162 3.5 Zero-Shot Generalizable Robotic Manipulation

163 As described in Section 3.3, we denote initial tracked points and features as s^0 and \mathbf{f}^0 . We estimate
 164 $\mathbf{s}_{\text{goal}} \in \mathbb{R}^{2 \times n_s}$ of goal image $\mathcal{I}_{\text{goal}}$ as follows:

$$\alpha_{ij} = \exp(\|\mathcal{W}_{\text{goal}}^f[\mathbf{u}_i] - \mathbf{f}_j^0\|_2),$$

$$w_{ij} = \frac{\exp(s\alpha_{ij})}{\sum_{i=1}^{H \times W} \exp(s\alpha_{ij})}, \quad (9)$$

165 then we have $\mathbf{s}_{\text{goal},j} = \sum_{i=1}^{H \times W} w_{ij} \mathbf{u}_i$, where $\mathcal{W}_{\text{goal}}^f$ is the feature volume extracted from $\mathcal{I}_{\text{goal}}$ using
 166 DINOv2. s is the hyperparameter to determine whether the heatmap w_{ij} is more smooth or concen-
 167 trating. Although Eq. 9 only shows a single instance case, it could be naturally extended to multiple
 168 instances by using instance mask information.

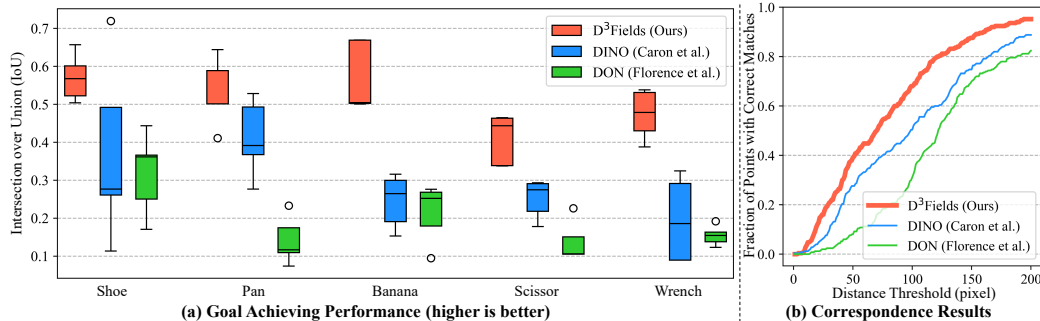


Figure 5: **Quantitative Evaluation.** We perform real-world quantitative evaluations by measuring final goal-achieving performance and keypoints correspondence accuracy. (a) We use IoU to measure goal-achieving performance. Results indicate that our method aligns with the goal configurations much better than DON and DINO across various object categories and scenarios. (b) We measure the keypoints correspondence accuracy according to the fraction of points with accurate matches, with correct matches determined by a distance threshold. Our method is consistently better at aligning with the goal image, regardless of the chosen threshold.

169 However, \mathbf{s}_{goal} is in the image space, while \mathbf{s}^t is in the 3D space. We bridge this gap by introducing
 170 a reference camera with approximate intrinsic and extrinsic parameters \mathbf{K}' and \mathbf{T}' . Instead of ren-
 171 dering images in the reference view, we focus on projecting 3D keypoints into 2D images and define
 172 the task cost function in image space as follows:

$$c(\mathbf{s}^t, \mathbf{s}_{\text{goal}}) = \|\pi(\mathbf{K}'(\mathbf{R}'\mathbf{s}^t + \mathbf{t}')) - \mathbf{s}_{\text{goal}}\|_2^2. \quad (10)$$

173 4 Experiments

174 In this section, we evaluate our representation across various manipulation tasks with varying goal
 175 image styles, instances, and contexts. We visualize D³Fields and showcase tracking results in Sec-
 176 tion 4.2. Then, we highlight our framework’s zero-shot generalizability in both real-world and
 177 simulated tasks in Section 4.3. Finally, a quantitative comparison with baselines in Section 4.4
 178 underscores our framework’s generalization and manipulation precision.

179 4.1 Experiment Setup

180 In the real world, we employ four OAK-PRO D cameras to gather RGBD observations and use the
 181 Kinova® Gen3 for action execution. In simulation, we utilize OmniGibson and deploy Fetch for
 182 mobile manipulation tasks [60]. Our evaluations span a variety of tasks, including organizing shoes,
 183 collecting debris, tidying the office table, arranging utensils, and more.

184 We implement the baseline methods using Dense Object Nets (DON) and DINO for feature extrac-
 185 tion [61, 54]. We quantitatively evaluate these methods on five object classes for single-instance
 186 manipulation tasks in the real world. The results and analysis are presented in Section 4.4.

187 4.2 Descriptor Fields Visualization and Keypoints Tracking

188 D³Fields provide a good 3D semantic representation, as shown in Fig. 3(a). We first visualize the
 189 mask fields by coloring 3D points according to their most likely instance, and our visualization
 190 shows a clear 3D instance segmentation. Additionally, we map the semantic features to RGB space
 191 using PCA, as with DINOv2 [17]. Visualization of the descriptor fields reveals that D³Fields retain
 192 a dense semantic understanding of objects. In the provided shoe example, even though various shoes
 193 have distinct appearances and poses, they exhibit similar color patterns: shoe heels are represented
 194 in green, and shoe toes in red. We observed similar patterns when evaluating the model on mugs
 195 and forks.

196 As discussed before, D³Fields can also capture scene dynamics. We evaluate it by tracking the
 197 object keypoints. We show two examples of 3D keypoint tracking in Fig. 3(b). In the first example,

198 a shoe is pushed and then flipped. Although only a portion of the shoe is visible from the view, our
199 framework tracks it reliably. In another example, a shoe is lifted and then set down. Despite parts of
200 the shoe being out of the camera’s view, we can robustly track it in 3D.

201 4.3 Zero-Shot Generalizable Manipulation

202 We conduct a qualitative evaluation of D³Fields in common household robotic manipulation tasks
203 in a zero-shot manner, with partial results displayed in Fig. 1 and Fig. 4. The following capabilities
204 of our framework are observed:

205 **Generalization to AI-Generated Goal Images.** In Fig. 1, the goal image, rendered in a Van Gogh
206 style, depicts shoes distinct from those in the workspace. Since D³Fields encode semantic informa-
207 tion, capturing shoes with varied appearances under similar descriptors, our framework can manip-
208 ulate shoes based on AI-generated goal images.

209 **Compositional Goal Images and 3D Manipulation.** Using the office desk organization example
210 in Fig. 1, the robot first arranges the mouse and pen according to the goal image. It then repositions
211 the mug from the box to the mug pad, referencing a goal image of the upright mug.

212 **Generalization across Instances and Materials.** Granular objects, unlike rigid ones, have more
213 complex dynamics. Our framework effectively handles these materials, as shown in the debris col-
214 lection in Fig. 1. Fig. 4 further showcases our framework’s instance-level generalization, where the
215 goal image displays instances different from the workspace.

216 **Generalization across Simulation and Real World.** We evaluated our framework on household
217 tasks in the simulator, as shown in the utensil organization and mug organization examples in Fig. 4.
218 Given goal images taken from the real world, our framework can also manipulate objects to the goal
219 configurations. Our framework demonstrates generalization capabilities between simulation and the
220 real world.

221 4.4 Quantitative Comparisons with Baselines

222 In Fig. 5(a), we measure performance using the IoU between the goal image mask and the final
223 state mask after manipulation, with higher values indicating better alignment. Evaluating across five
224 object classes, our method consistently outperforms the baselines, underscoring its generalization
225 and manipulation accuracy. While DINO struggles with distinguishing object components, leading
226 to imprecise results, it still works better than DON. Although DON performs well on familiar object
227 classes and configurations, it lacks generalization in novel scenarios.

228 In Fig. 5(b), we present the correspondence results. We manually label corresponding keypoints
229 on both the goal image and the final manipulation result to evaluate the correspondence accuracy.
230 We calculate the fraction of accurately matched points based on a distance threshold. Our method
231 consistently outperforms the baselines, regardless of the threshold. DINO ranks second, while DON
232 lags behind. Consistent with Fig. 5(a), our method excels in generalization and accuracy, DINO is
233 broadly applicable but less precise, and DON struggles with generalization.

234 5 Conclusion

235 In this work, we introduce D³Fields, which implicitly encode 3D semantic features and 3D instance
236 masks, and model the underlying dynamics. Our emphasis is on zero-shot generalizable robotic
237 manipulation tasks specified by 2D goal images of varying styles, contexts, and instances. Our
238 framework excels in executing a diverse array of household manipulation tasks in both simulated
239 and real-world scenarios. Its performance greatly surpasses baseline methods such as Dense Object
240 Nets and DINO in terms of generalization capabilities and manipulation accuracy.

241 **References**

- 242 [1] L. Manuelli, Y. Li, P. Florence, and R. Tedrake. Keypoints into the future: Self-supervised
243 correspondence in model-based reinforcement learning. In *Conference on Robot Learning*
244 (*CoRL*), 2020.
- 245 [2] Y. Li, S. Li, V. Sitzmann, P. Agrawal, and A. Torralba. 3d neural scene representations for
246 visuomotor control. *arXiv preprint arXiv:2107.04004*, 2021.
- 247 [3] H. Shi, H. Xu, Z. Huang, Y. Li, and J. Wu. Robocraft: Learning to see, simulate, and shape
248 elasto-plastic objects with graph networks. *arXiv preprint arXiv:2205.02909*, 2022.
- 249 [4] Y. Ze, N. Hansen, Y. Chen, M. Jain, and X. Wang. Visual reinforcement learning with self-
250 supervised 3d representations. *IEEE Robotics and Automation Letters (RA-L)*, 2023.
- 251 [5] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent
252 dynamics for planning from pixels. In *International conference on machine learning*, pages
253 2555–2565. PMLR, 2019.
- 254 [6] W. Yan, A. Vangipuram, P. Abbeel, and L. Pinto. Learning predictive representations for de-
255 formable objects using contrastive estimation. In J. Kober, F. Ramos, and C. Tomlin, editors,
256 *Proceedings of the 2020 Conference on Robot Learning*, volume 155 of *Proceedings of Ma-*
257 *chine Learning Research*, pages 564–574. PMLR, 16–18 Nov 2021.
- 258 [7] Y. Wang, Y. Li, K. Driggs-Campbell, L. Fei-Fei, and J. Wu. Dynamic-Resolution Model
259 Learning for Object Pile Manipulation. In *Proceedings of Robotics: Science and Systems*,
260 Daegu, Republic of Korea, July 2023. doi:10.15607/RSS.2023.XIX.047.
- 261 [8] M. Minderer, C. Sun, R. Villegas, F. Cole, K. P. Murphy, and H. Lee. Unsupervised learning
262 of object structure and dynamics from videos. *Advances in Neural Information Processing*
263 *Systems*, 32, 2019.
- 264 [9] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield. Deep object pose
265 estimation for semantic robotic grasping of household objects. In *Conference on Robot Learn-*
266 *ing*, pages 306–316. PMLR, 2018.
- 267 [10] S. Tyree, J. Tremblay, T. To, J. Cheng, T. Mosier, J. Smith, and S. Birchfield. 6-dof pose esti-
268 mation of household objects for robotic manipulation: An accessible dataset and benchmark.
269 In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages
270 13081–13088. IEEE, 2022.
- 271 [11] K. M. Jatavallabhula, A. Kuwajerwala, Q. Gu, M. Omama, G. Iyer, S. Saryazdi, T. Chen,
272 A. Maalouf, S. Li, N. V. Keetha, A. Tewari, J. Tenenbaum, C. de Melo, M. Krishna, L. Paull,
273 F. Shkurti, and A. Torralba. ConceptFusion: Open-set multimodal 3D mapping. In *Proceedings*
274 *of Robotics: Science and Systems*, Daegu, Republic of Korea, July 2023. doi:10.15607/RSS.
275 2023.XIX.066.
- 276 [12] K. Mazur, E. Sucar, and A. J. Davison. Feature-realistic neural fusion for real-time, open set
277 scene understanding. In *2023 IEEE International Conference on Robotics and Automation*
278 (*ICRA*), pages 8201–8207. IEEE, 2023.
- 279 [13] W. Liu, Y. Du, T. Hermans, S. Chernova, and C. Paxton. Structdiffusion: Language-guided
280 creation of physically-valid structures using unseen objects. In *RSS 2023*, 2023.
- 281 [14] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, et al. Grounding
282 dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint*
283 *arXiv:2303.05499*, 2023.
- 284 [15] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead,
285 A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick. Segment anything. *arXiv:2304.02643*, 2023.

- 286 [16] H. K. Cheng and A. G. Schwing. XMem: Long-term video object segmentation with an
287 atkinson-shiffrin memory model. In *ECCV*, 2022.
- 288 [17] M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez,
289 D. Haziza, F. Massa, A. El-Nouby, R. Howes, P.-Y. Huang, H. Xu, V. Sharma, S.-W. Li,
290 W. Galuba, M. Rabbat, M. Assran, N. Ballas, G. Synnaeve, I. Misra, H. Jegou, J. Mairal,
291 P. Labatut, A. Joulin, and P. Bojanowski. Dinov2: Learning robust visual features without
292 supervision, 2023.
- 293 [18] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang,
294 R. Julian, et al. Do as i can, not as i say: Grounding language in robotic affordances. In
295 *Conference on Robot Learning*, pages 287–318. PMLR, 2023.
- 296 [19] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch,
297 Y. Chebotar, et al. Inner monologue: Embodied reasoning through planning with language
298 models. In *Conference on Robot Learning*, pages 1769–1782. PMLR, 2023.
- 299 [20] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code
300 as policies: Language model programs for embodied control. In *2023 IEEE International
301 Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023.
- 302 [21] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei. Voxposer: Composable 3d value
303 maps for robotic manipulation with language models. In *7th Annual Conference on Robot
304 Learning*, 2023.
- 305 [22] Y. Zhu, Z. Jiang, P. Stone, and Y. Zhu. Learning generalizable manipulation policies with
306 object-centric 3d representations. In *7th Annual Conference on Robot Learning*, 2023.
- 307 [23] K. Mülling, J. Kober, O. Kroemer, and J. Peters. Learning to select and generalize striking
308 movements in robot table tennis. *The International Journal of Robotics Research*, 32(3):263–
309 279, 2013.
- 310 [24] Y. Duan, M. Andrychowicz, B. Stadie, O. Jonathan Ho, J. Schneider, I. Sutskever, P. Abbeel,
311 and W. Zaremba. One-shot imitation learning. *Advances in neural information processing
312 systems*, 30, 2017.
- 313 [25] Y. Li, J. Wu, R. Tedrake, J. B. Tenenbaum, and A. Torralba. Learning particle dynamics for
314 manipulating rigid bodies, deformable objects, and fluids. In *ICLR*, 2019.
- 315 [26] W. Wang, A. S. Morgan, A. M. Dollar, and G. D. Hager. Dynamical scene representation
316 and control with keypoint-conditioned neural radiance field. In *2022 IEEE 18th International
317 Conference on Automation Science and Engineering (CASE)*, pages 1138–1143. IEEE, 2022.
- 318 [27] W. Gao and R. Tedrake. kpm 2.0: Feedback control for category-level robotic manipulation.
319 *IEEE Robotics and Automation Letters*, 6(2):2962–2969, 2021.
- 320 [28] L. Manuelli, W. Gao, P. Florence, and R. Tedrake. kpm: Keypoint affordances for category-
321 level robotic manipulation. In *The International Symposium of Robotics Research*, pages 132–
322 157. Springer, 2019.
- 323 [29] W. Gao and R. Tedrake. kpm-sc: Generalizable manipulation planning using keypoint af-
324 fordance and shape completion. In *2021 IEEE International Conference on Robotics and
325 Automation (ICRA)*, pages 6527–6533. IEEE, 2021.
- 326 [30] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent
327 imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 328 [31] X. Lin, C. Qi, Y. Zhang, Y. Li, Z. Huang, K. Fragkiadaki, C. Gan, and D. Held. Planning
329 with spatial-temporal abstraction from point clouds for deformable object manipulation. In *6th
330 Annual Conference on Robot Learning*, 2022.

- 331 [32] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual repre-
332 sentation for robot manipulation. In *Conference on Robot Learning*, pages 892–909. PMLR,
333 2023.
- 334 [33] T. Xiao, I. Radosavovic, T. Darrell, and J. Malik. Masked visual pre-training for motor control.
335 *arXiv:2203.06173*, 2022.
- 336 [34] I. Radosavovic, T. Xiao, S. James, P. Abbeel, J. Malik, and T. Darrell. Real-world robot
337 learning with masked visual pre-training. *CoRL*, 2022.
- 338 [35] Z. Mandi, H. Bharadhwaj, V. Moens, S. Song, A. Rajeswaran, and V. Kumar. Cacti:
339 A framework for scalable multi-task multi-scene visual imitation learning. *arXiv preprint*
340 *arXiv:2212.05711*, 2022.
- 341 [36] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, B. Zitkovich,
342 F. Xia, C. Finn, and K. Hausman. Open-world object manipulation using pre-trained vision-
343 language models, 2023.
- 344 [37] Y. Yoon, G. N. DeSouza, and A. C. Kak. Real-time tracking and pose estimation for industrial
345 objects using geometric features. In *2003 IEEE International conference on robotics and*
346 *automation (cat. no. 03CH37422)*, volume 3, pages 3473–3478. IEEE, 2003.
- 347 [38] M. Zhu, K. G. Derpanis, Y. Yang, S. Brahmabhatt, M. Zhang, C. Phillips, M. Lecce, and
348 K. Daniilidis. Single image 3d object detection and pose estimation for grasping. In *2014*
349 *IEEE International Conference on Robotics and Automation (ICRA)*, pages 3936–3943. IEEE,
350 2014.
- 351 [39] L. Zhu, A. Mousavian, Y. Xiang, H. Mazhar, J. van Eenbergen, S. Debnath, and D. Fox. Rgb-
352 d local implicit function for depth completion of transparent objects. In *Proceedings of the*
353 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4649–4658, 2021.
- 354 [40] J. Ichnowski, Y. Avigal, J. Kerr, and K. Goldberg. Dex-nerf: Using a neural radiance field to
355 grasp transparent objects. In *5th Annual Conference on Robot Learning*, 2021.
- 356 [41] Y. Wi, P. Florence, A. Zeng, and N. Fazeli. VirDo: Visio-tactile implicit representations of
357 deformable objects. In *2022 International Conference on Robotics and Automation (ICRA)*,
358 pages 3583–3590. IEEE, 2022.
- 359 [42] A. Simeonov, Y. Du, A. Tagliasacchi, J. B. Tenenbaum, A. Rodriguez, P. Agrawal, and V. Sitz-
360 mann. Neural descriptor fields: Se (3)-equivariant object representations for manipulation. In
361 *2022 International Conference on Robotics and Automation (ICRA)*, pages 6394–6400. IEEE,
362 2022.
- 363 [43] D. Driess, J.-S. Ha, M. Toussaint, and R. Tedrake. Learning models as functionals of signed-
364 distance fields for manipulation planning. In *Conference on Robot Learning*, pages 245–255.
365 PMLR, 2022.
- 366 [44] Z. Jiang, Y. Zhu, M. Svetlik, K. Fang, and Y. Zhu. Synergies Between Affordance and Geome-
367 try: 6-DoF Grasp Detection via Implicit Representations. In *Proceedings of Robotics: Science*
368 *and Systems*, Virtual, July 2021. doi:10.15607/RSS.2021.XVII.024.
- 369 [45] T. Weng, D. Held, F. Meier, and M. Mukadam. Neural grasp distance fields for robot manipu-
370 lation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- 371 [46] D. Driess, I. Schubert, P. Florence, Y. Li, and M. Toussaint. Reinforcement learning with neural
372 radiance fields. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 373 [47] D. Shim, S. Lee, and H. J. Kim. SNeRL: Semantic-aware neural radiance fields for reinforc-
374 e-ment learning. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett,
375 editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202
376 of *Proceedings of Machine Learning Research*, pages 31489–31503. PMLR, 23–29 Jul 2023.

- 377 [48] L. Yen-Chen, P. Florence, J. T. Barron, T.-Y. Lin, A. Rodriguez, and P. Isola. NeRF-
378 Supervision: Learning dense object descriptors from neural radiance fields. In *IEEE Con-*
379 *ference on Robotics and Automation (ICRA)*, 2022.
- 380 [49] Z. Tang, B. Sundaralingam, J. Tremblay, B. Wen, Y. Yuan, S. Tyree, C. Loop, A. Schwing,
381 and S. Birchfield. RGB-only reconstruction of tabletop scenes for collision-free manipulator
382 control. In *ICRA*, 2023.
- 383 [50] A. Zhou, M. J. Kim, L. Wang, P. Florence, and C. Finn. Nerf in the palm of your hand: Cor-
384 rective augmentation for robotics via novel-view synthesis. In *Proceedings of the IEEE/CVF*
385 *Conference on Computer Vision and Pattern Recognition*, pages 17907–17917, 2023.
- 386 [51] N. M. (Mahi)Shafiullah, C. Paxton, L. Pinto, S. Chintala, and A. Szlam. CLIP-Fields: Weakly
387 Supervised Semantic Fields for Robotic Memory. In *Proceedings of Robotics: Science and*
388 *Systems*, Daegu, Republic of Korea, July 2023. doi:10.15607/RSS.2023.XIX.074.
- 389 [52] Y. Wi, A. Zeng, P. Florence, and N. Fazeli. VirDo++: Real-world, visuo-tactile dynamics and
390 perception of deformable objects. *arXiv preprint arXiv:2210.03701*, 2022.
- 391 [53] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
392 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from
393 natural language supervision. In M. Meila and T. Zhang, editors, *Proceedings of the 38th Inter-*
394 *national Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning*
395 *Research*, pages 8748–8763. PMLR, 18–24 Jul 2021.
- 396 [54] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging
397 properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF interna-*
398 *tional conference on computer vision*, pages 9650–9660, 2021.
- 399 [55] J. Kerr, C. M. Kim, K. Goldberg, A. Kanazawa, and M. Tancik. Lerf: Language embedded
400 radiance fields. In *International Conference on Computer Vision (ICCV)*, 2023.
- 401 [56] S. Sharma, A. Rashid, C. M. Kim, J. Kerr, L. Y. Chen, A. Kanazawa, and K. Goldberg. Lan-
402 guage embedded radiance fields for zero-shot task-oriented grasping. In *7th Annual Conference*
403 *on Robot Learning*, 2023.
- 404 [57] W. Shen, G. Yang, A. Yu, J. Wong, L. P. Kaelbling, and P. Isola. Distilled feature fields enable
405 few-shot manipulation. In *7th Annual Conference on Robot Learning*, 2023.
- 406 [58] Y. Ze, G. Yan, Y.-H. Wu, A. Macaluso, Y. Ge, J. Ye, N. Hansen, L. E. Li, and X. Wang. Multi-
407 task real robot learning with generalizable neural feature fields. In *7th Annual Conference on*
408 *Robot Learning*, 2023.
- 409 [59] Y. Li, T. Lin, K. Yi, D. Bear, D. L. Yamins, J. Wu, J. B. Tenenbaum, and A. Torralba. Visual
410 grounding of learned physical models. In *International Conference on Machine Learning*,
411 2020.
- 412 [60] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine,
413 M. Lingelbach, J. Sun, M. Anvari, M. Hwang, M. Sharma, A. Aydin, D. Bansal, S. Hunter,
414 K.-Y. Kim, A. Lou, C. R. Matthews, I. Villa-Renteria, J. H. Tang, C. Tang, F. Xia, S. Savarese,
415 H. Gweon, K. Liu, J. Wu, and L. Fei-Fei. BEHAVIOR-1k: A benchmark for embodied AI
416 with 1,000 everyday activities and realistic simulation. In *6th Annual Conference on Robot*
417 *Learning*, 2022.
- 418 [61] P. R. Florence, L. Manuelli, and R. Tedrake. Dense object nets: Learning dense visual ob-
419 ject descriptors by and for robotic manipulation. In A. Billard, A. Dragan, J. Peters, and
420 J. Morimoto, editors, *Proceedings of The 2nd Conference on Robot Learning*, volume 87 of
421 *Proceedings of Machine Learning Research*, pages 373–385. PMLR, 29–31 Oct 2018.