

# NeRF in the Palm of Your Hand: Corrective Augmentation for Robotics via Novel-View Synthesis

Allan Zhou\*  
Stanford

Moo Jin Kim\*  
Stanford

Lirui Wang  
MIT CSAIL

Pete Florence  
Google

Chelsea Finn  
Stanford

## Abstract

Expert demonstrations are a rich source of supervision for training visual robotic manipulation policies, but imitation learning methods often require either a large number of demonstrations or expensive online expert supervision to learn reactive closed-loop behaviors. In this work, we introduce SPARTN (Synthetic Perturbations for Augmenting Robot Trajectories via NeRF): a fully-offline data augmentation scheme for improving robot policies that use eye-in-hand cameras. Our approach leverages neural radiance fields (NeRFs) to synthetically inject corrective noise into visual demonstrations, using NeRFs to generate perturbed viewpoints while simultaneously calculating the corrective actions. This requires no additional expert supervision or environment interaction, and distills the geometric information in NeRFs into a real-time reactive RGB-only policy. In a simulated 6-DoF visual grasping benchmark, SPARTN improves success rates by  $2.8\times$  over imitation learning without the corrective augmentations and even outperforms some methods that use online supervision. It additionally closes the gap between RGB-only and RGB-D success rates, eliminating the previous need for depth sensors. In real-world 6-DoF robotic grasping experiments from limited human demonstrations, our method improves absolute success rates by 22.5% on average, including objects that are traditionally challenging for depth-based methods. See video results at <https://bland.website/spartn>.

## 1. Introduction

Object grasping is a central problem in vision-based control and is fundamental to many robotic manipulation problems. While there has been significant progress in top-down bin picking settings [21, 34], 6-DoF grasping of arbitrary objects amidst clutter remains an open problem, and is especially challenging for shiny or reflective objects that are not visible to depth cameras. For example, the task of grasping a wine glass from the stem shown in Figure 1 re-

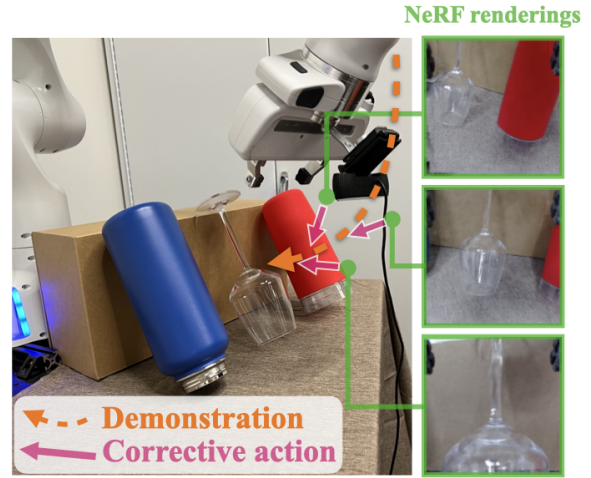


Figure 1. SPARTN is an offline data augmentation method for behavior cloning eye-in-hand visual policies. It simulates recovery in a **demonstration** by using NeRFs to render high-fidelity **observations** (right) from noisy states, then generates **corrective action** labels.

quires precise 6-DoF control (using full 3D translation and 3D rotation of the gripper) and closed-loop perception of a transparent object. Traditional 6-DoF grasping pipelines [8, 57] synthesize only one grasp pose and use a motion planner to generate a collision-free trajectory to reach the grasp [38, 40, 54, 60]. However, the use of open-loop trajectory execution prevents the system from using perceptual feedback for reactive, precise grasping behavior. In this paper, we study how to learn closed-loop policies for 6-DoF object grasping from RGB images, which can be trained with imitation or reinforcement learning methods [58].

Imitation learning from expert demonstrations is a simple and promising approach to this problem, but is known to suffer from compounding errors [46]. As a result, complex vision-based tasks can require online expert supervision [19, 46] or environment interaction [13, 44], both of which are expensive and time-consuming to collect. On the other hand, offline “feedback augmentation” methods [14, 22] can be effective at combating compounding errors, but are severely limited in scope and thus far have not been ap-

\*Equal contribution. Correspondence to [ayz@stanford.edu](mailto:ayz@stanford.edu).

plied to visual observations. Other recent works have found that using eye-in-hand cameras mounted on a robot’s wrist can significantly improve the performance of visuomotor policies trained with imitation learning [17, 20, 35], but still do not address the underlying issue of compounding errors. We develop an approach that helps address compounding errors to improve vision-based policies, while building on the success of eye-in-hand cameras.

To improve imitation learning for quasi-static tasks like grasping, we propose a simple yet effective offline data augmentation technique. For an eye-in-hand camera, the images in each demonstration trajectory form a collection of views of the demonstration scene, which we use to train neural radiance fields (NeRFs) [37] of each scene. Then, we can augment the demonstration data with corrective feedback by injecting noise into the camera poses along the demonstration and using the demonstration’s NeRF to render observations from the new camera pose. Because the camera to end-effector transform is known, we can compute corrective action labels for the newly rendered observations by considering the action that would return the gripper to the expert trajectory. The augmented data can be combined with the original demonstrations to train a reactive, real-time policy. Since the NeRFs are trained on the original demonstrations, this method effectively “distills” the 3D information from each NeRF into the policy.

The main contribution of this work is a NeRF-based data augmentation technique, called SPARTN (Synthetic Perturbations for Augmenting Robot Trajectories via NeRF), that improves behavior cloning for eye-in-hand visual grasping policies. By leveraging view-synthesis methods like NeRF, SPARTN extends the idea of corrective feedback augmentation to the visual domain. The resulting approach can produce (i) reactive, (ii) real-time, and (iii) RGB-only policies for 6-DoF grasping. The data augmentation is fully offline and does not require additional effort from expert demonstrators nor online environment interactions. We evaluate SPARTN on 6-DoF robotic grasping tasks both in simulation and in the real world. On a previously-proposed simulated 6-DoF grasping benchmark [58], the augmentation from SPARTN improves grasp success rates by  $2.8\times$  compared to training without SPARTN, and even outperforms some methods that use expensive online supervision. On eight challenging real-world grasping tasks with a Franka Emika Panda robot, SPARTN improves the absolute average success rate by 22.5%.

## 2. Related Work

**Robotic Grasping.** Grasping is a long-studied topic in robotics [50]; see the multitude of survey articles for a complete review [2, 4, 26]. Most data-driven grasping systems focus on learning how to predict some parameterized “grasp” (whether a full 6-DoF pose, 2-DoF table-top po-

Method	No Depth	Full 6-DoF	Closed-Loop
[34, 45]			
[53]		✓	✓
[18, 24, 40, 54]		✓	
[21, 32]	✓		✓
[43]	✓		
[39]			✓
[56, 58]		✓	✓
SPARTN (ours)	✓	✓	✓

Table 1. Comparison of our approach with related grasping work. SPARTN is the only approach to learn closed-loop 6-DoF grasping policies from only RGB inputs.

sition, etc.), and leave intermediate motion generation to be open-loop, handled either through motion planners or simple heuristics, e.g. [12, 34, 42, 43, 52, 57]. Other works have trained *closed-loop* grasping policies [21, 39, 53, 58], and bring all the benefits of closed-loop policies: including for example, the ability to avoid obstacles, to perform precise grasping without precise calibration, and to react to dynamic objects. Additionally, grasping policies are often designed for top-down (2- or 3-DoF) grasping [21, 31, 43], while 6-DoF grasping typically requires depth or 3D information [39, 42, 53, 58]. In contrast, our method trains a reactive 6-DoF grasping policy with only RGB data. See Table 1 for a summary of the assumptions of the most related grasping works.

**Imitation learning and data augmentation.** Behavior cloning is known to struggle with covariate shift: small errors cause imitation policies to fall slightly off of the data distribution and it is then difficult to correct the mistake back onto the data manifold. DAgger [46] and its variants [16, 23, 36] mitigate this issue by obtaining expert corrections throughout training. Alternatively, DART [29] injects noise during expert demonstration collection, which is especially effective with algorithmic experts but interferes with human demonstration collection. Previous works [14, 22] have injected noise into the low-dimensional system state after data collection (in a fully offline manner), but the visual observations are left out, limiting the helpfulness of noise injection. Our method can be seen as a visual, fully-offline version of noise injection that does not require perturbing the expert during demonstration collection, using NeRF to synthetically render perturbed states post-hoc. Unlike standard image augmentations for policy learning [30, 61], our method uses NeRF to learn a 3D model of the demonstration scene, which enables us to generate high-fidelity novel views for data augmentation. In addition, while standard image augmentation approaches do not modify the action labels, we leverage hand-eye coordination to calculate corrective actions for augmented observations.

**NeRF for Robotics.** A number of recent works have in-

vestigated applications of NeRF and related methods in robotics, including localization [63], navigation [1], dynamics modeling [7, 9, 33], reinforcement learning [10], and data generation for other learning-based methods [18, 62]. NeRF-Supervision [62], for example, generates pixel-level correspondence to learn dense object descriptors, which are useful for manipulation tasks. For grasping, various methods have leveraged NeRF [3, 18, 24] for open-loop grasp synthesis. In contrast, our method uses NeRF offline to augment data for grasping and distills a reactive, real-time, RGB-only, closed-loop policy.

### 3. Methodology

We now introduce SPARTN, which augments an eye-in-hand robot demonstration dataset using NeRF. We first review preliminaries, then describe a method overview, followed by details of training NeRFs and augmenting corrective behavior.

#### 3.1. Preliminaries

**Imitation learning.** In imitation learning, we assume access to a dataset of  $N$  expert trajectories  $\mathcal{D} = \{\tau\}_{i=1}^N$ , where each trajectory consists of a sequence of state-action pairs,  $\tau = \{(s_k, a_k)\}_{k=0}^K$ . In purely *offline* imitation learning, there exists no other primary data-collection assumptions other than this dataset  $\mathcal{D}$ , i.e. no reward labels or online interactions. A standard method for this offline setting is behavior cloning (BC), which trains a policy  $\pi_\theta$  to mimic the expert via supervised learning, by minimizing the objective  $\mathcal{L}(\theta) = \mathbb{E}_{(s,a) \sim \mathcal{D}}[\ell(\pi_\theta(s), a)]$ , where  $\ell$  is some loss function in the action space.

**NeRF.** Our method uses novel-view synthesis as a building block, for which we use Neural Radiance Fields [37]. For each scene, NeRF performs novel-view synthesis by training a scene-specific neural radiance field  $F_\Theta$  from a “training set” of posed images  $\{(I_k, T_k)\}$ , where each  $I_k \in \mathbb{R}^{w \times h \times 3}$ ,  $T_k \in SE(3)$ . After  $F_\Theta$  is trained, through volume rendering NeRF can render new views of a scene from any requested pose, which can be summarized as  $I = \text{NeRF-Render}(T; F_\Theta)$  – in particular, this works best “near” the training set of poses. Since we train many NeRFs (one per demonstration), we use an accelerated implementation of NeRF (Instant-NeRF [41]).

**Corrective Noise Augmentation.** A simple method which has been shown to improve the robustness of behavior-cloned policies is to perform corrective noise augmentation [14, 22]. The idea is to take a nominal trajectory  $\tau = \{(s_k, a_k)\}_{k=0}^K$  and create a noise-distributed corrective trajectory  $\tilde{\tau} = \{(\tilde{s}_k, \tilde{a}_k)\}_{k=0}^K$ , where each sampled state  $\tilde{s}_k$  is a noisy version of the measured state, i.e.  $\tilde{s}_k \sim s_k + \epsilon$  and  $\epsilon$  is sampled noise. To perform corrective feedback augmentation,  $\tilde{a}_k$  is chosen such that it will return the state to the nominal trajectory: given the true inverse dynamics  $f^{-1}$

---

#### Algorithm 1 SPARTN (Simplified overview)

---

**Input:**  $\mathcal{D} = \{\tau\}_{i=1}^N$  - expert demonstrations  
**Output:** Augmented transition dataset  $\tilde{\mathcal{D}}$

- 1:  $\tilde{\mathcal{D}} \leftarrow \{\}$
- 2: **for**  $\tau = \{(I_k, T_k, a_k)\}_{k=1}^K \in \mathcal{D}$  **do**
- 3:  $F_\Theta^\tau \leftarrow \text{NeRF}(\{(I_k, T_k)\}_{k=1}^K)$  # train NeRF
- 4: **for**  $k \in (0, 1, \dots, K)$  **do**
- 5:     **for**  $i = 1 : N_{\text{aug}}$  **do**
- 6:          $\epsilon \leftarrow \text{NoiseDist}(SE(3))$
- 7:          $\tilde{T}_k \leftarrow \text{Perturb}(T_k, \epsilon)$  # perturb pose
- 8:          $\tilde{a}_k \leftarrow \text{CorrectAction}(T_k, \tilde{T}_k, a_k)$
- 9:          $\tilde{I}_k \leftarrow \text{NeRF-Render}(\tilde{T}_k, F_\Theta^\tau)$
- 10:          $\tilde{\mathcal{D}} \leftarrow \tilde{\mathcal{D}} \cup \{(\tilde{I}_k, \tilde{T}_k, \tilde{a}_k)\}$

---

of the environment, then  $\tilde{a}_k = f^{-1}(\tilde{s}_k, s_{k+1})$ , or compute the commanded control that can be tracked with a stabilizing controller to  $s_{k+1}$ . An easy way to parameterize this is by choosing the action space of the learned policy to be the input to a stabilizing controller, as in [14, 22]. Note that it is also common in the imitation and reinforcement learning literature to apply noise to inputs, which is typically interpreted as a way to regularize the policy [27, 30]. Meanwhile, in addition to potential regularization effects, the *corrective* noise augmentation has been interpreted to specifically reduce compounding errors [14, 22], but of course has limits if the scale of perturbations is too large or in highly non-smooth dynamics regimes. A critical limitation of prior works using corrective noise augmentation is that they have not been applied to visual observations.

#### 3.2. Overview: Visual Corrective Augmentation

Consider the case where an agent receives partial visual observations  $I$  instead of the full state of the world, and also receives direct proprioceptive state measurements,  $s^{\text{robot}}$ , as is commonly the case in robotics. In general, the corrective augmentation of Sec. 3.1 requires obtaining the visual observation  $\tilde{I}_k$  for each noisy state  $\tilde{s}_k$ , which could be expensive or impossible in the actual environment.

An insight of this work is that for *eye-in-hand* robot policies (where the visual observation comes from a camera mounted on the wrist) in static scenes, we can readily generate novel visual observations using novel-view synthesis (i.e., NeRF) without further interactions in the environment. In this setting, a primary subset of the *observations* are posed images  $(I, T)$ , together with some *actions*  $a$ .

The key intuition of our method can be grasped by considering how to perform visually corrective augmentation via NeRF, as illustrated in Figure 2. Noisy states and corrective actions  $(\tilde{T}_k, \tilde{a}_k)$  can be generated for a trajectory of posed observations and actions  $\tau = \{(I_k, T_k, a_k)\}_{k=0}^K$  (Sec. 3.1). The key pre-processing step is to train a trajectory-specific NeRF  $F_\Theta^\tau$  for each demonstration  $\tau$  (Sec. 3.3). These trajectory-specific NeRFs enable us to

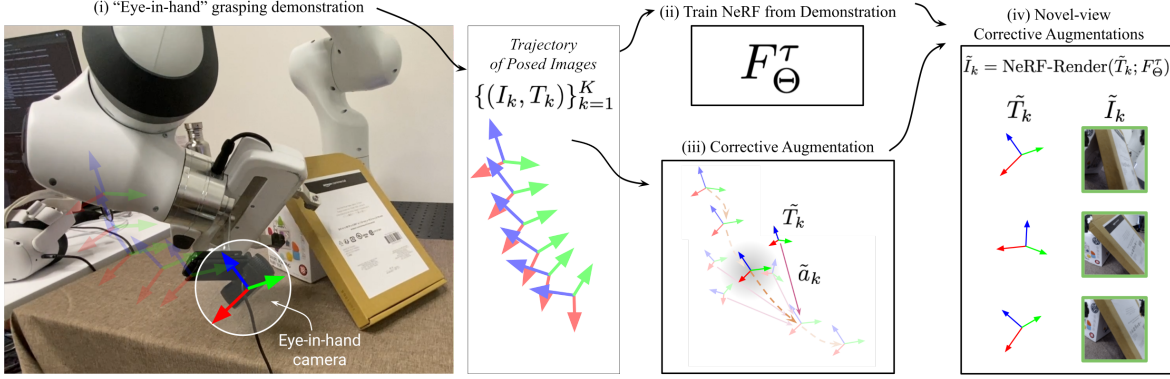


Figure 2. An illustration of how SPARTN creates augmentations from an original demonstration (in reality, this process is repeated for every available demonstration). **(i)**: The eye-in-hand demonstration contains posed images  $\{(I_k, T_k)\}_{k=1}^K$ . **(ii)**: We train a neural radiance field (NeRF) of the demonstration scene on the posed images. **(iii)**: We sample perturbations around each pose to **simulate** noise in the demonstration, and calculate the corrective action (in magenta) that would stabilize the trajectory. **(iv)**: We use the NeRF to render observations for the perturbed poses. The end result is augmented image-action pairs for improving behavior cloning.

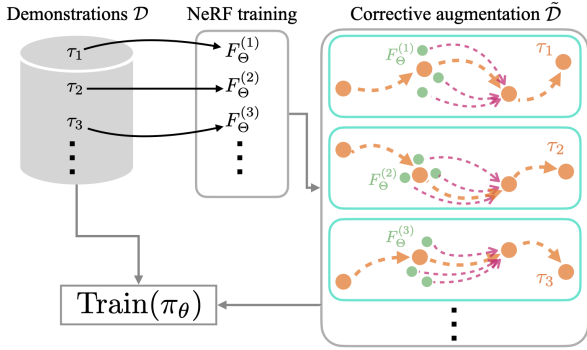


Figure 3. An overview of the SPARTN training process. A NeRF is trained for each of the original demonstrations in  $\mathcal{D}$ . We use these NeRFs to generate visual corrective augmentations for each demonstration and collect them in  $\tilde{\mathcal{D}}$ . The policy  $\pi_{\theta}$  can be trained on  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$  using standard behavior cloning methods.

render observations  $\tilde{I}_k$  for noisy states  $\tilde{T}_k$ , completing the augmentation process and resulting in visually corrective transitions  $(\tilde{I}_k, \tilde{T}_k, \tilde{a}_k)$  (Sec. 3.4). Algorithm 1 and Figure 3 overview this process.

### 3.3. Training NeRFs from Robot Demonstrations

SPARTN uses novel-view synthesis, in particular NeRF, to generate observations  $\tilde{I}_k$  for noisy states without environment interaction. We train a NeRF  $F_{\Theta}^T$  for each demonstration trajectory using the image observations  $(I_1, \dots, I_K) \in \tau$  as the training set of views. After training  $F_{\Theta}^T$ , we create observations  $\tilde{I}_k$  for perturbed robot states  $\tilde{T}_k$  by rendering the view from the perturbed camera pose using  $F_{\Theta}^T$ .

An important detail is that the end-effector reference frame used for control in demonstrations may differ from the reference frame for the camera itself, but through stan-

dard eye-in-hand calibration we can transform all visual observations used for training and augmenting the NeRFs into the camera frame. Given a transform  ${}^W T_k^{\text{from}}$  which transforms between two frames, we simply transform all NeRF-poses to the world frame:  ${}^W T_k^C = {}^W T_k^{EE} E T^C$ , where  $W$  is the world frame,  $E$  is the end-effector frame which changes at each step  $k$ , and  $C$  is the camera (NeRF) frame statically linked to the end-effector frame, and  ${}^E T^C$  is acquired through hand-eye calibration.

**COLMAP camera poses.** Real-world calibration error means that our camera-to-world transforms  $\{{}^W T_k^C\}_{k=1}^K$  are noisy and we obtain higher-quality NeRFs by using camera poses estimated by COLMAP [48, 49]. Up to noise and a scale factor  $\beta$ , the only difference from the world frame camera transforms is that COLMAP uses an arbitrary reference frame  $V \neq W$ . We denote the COLMAP outputs  $\{{}^V H_k^C\}_{k=1}^K$ , using  $H$  instead of  $T$  because of the difference in scale. We now introduce notation to separate the rotation and translation components of a transform:

$${}^a T^b := (\text{Rot} [{}^a T^b], \text{Trans} [{}^a T^b]) \quad (1)$$

Since we train the NeRFs on COLMAP’s camera poses, we must convert perturbed camera poses  ${}^W T_k^{\tilde{C}}$  to COLMAP’s frame in order to render the observations. In other words, we must call  $\text{NeRF-Render}({}^V H_k^{\tilde{C}}, F_{\Theta}^T)$ , where:

$${}^V H_k^{\tilde{C}} = \left( \text{Rot} [{}^V T_k^{\tilde{C}}], \beta \text{Trans} [{}^V T_k^{\tilde{C}}] \right) \quad (2)$$

$${}^V T_k^{\tilde{C}} = {}^V T^W W T_k^{\tilde{C}}. \quad (3)$$

Both  $\beta$  and  ${}^V T^W$  can be estimated from the pairs  $\{({}^W T_k^C, {}^V H_k^C)\}_{k=1}^K$ , as described in Appendix D.2.

**Static Scene Assumption.** An additional consideration for robotic manipulation is that the standard NeRF formulation

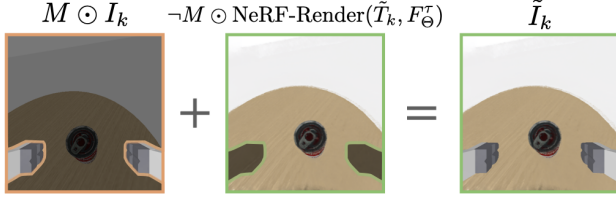


Figure 4. An illustration of how the gripper is inserted into the result of the NeRF rendering process. Gray regions indicate pixels being masked out by the binary gripper mask  $M \in \{0, 1\}^{w \times h}$ , which denotes the pixels where the gripper is located in all frames.

assumes a static scene, while manipulation tasks such as grasping will usually move objects in the scene. To address this, we apply the NeRF training and augmentation process to only the subsets of each demonstration trajectory where no robot-object interaction occurs, instead of all timesteps. For grasping, a simple and effective heuristic is to only apply SPARTN to the portions of each demonstration *before* the expert closes the gripper.

**Masking Out the Robot Gripper.** The robot’s gripper is often within the view of an eye-in-hand camera, which breaks the static scene assumption. To address this, we leverage that the gripper is in the same location in each image assuming the camera is rigidly mounted and the gripper is open. Further, since NeRF is trained per-ray, we can simply mask pixels from training images. We construct a single binary mask  $M \in \{0, 1\}^{w \times h}$ , where a 1 indicates gripper pixels to mask out. Figure 4 shows how we use the same mask to splice the gripper back into each NeRF rendering output,  $\tilde{I}_k \leftarrow \neg M \odot \tilde{I}_k + M \odot I_k$ , where  $\odot$  is element-wise multiplication broadcasted across the color channels.

**NeRF Quality.** Even when the training views from a demonstration are suboptimal for NeRF training, SPARTN benefits from the fact that our augmentation is local and the perturbations  $\varepsilon$  are typically small, so we only need the NeRF to generalize in a small region around the demonstration trajectory. As we will verify in the experiments, SPARTN can be effective even with a limited number of training views compared to other NeRF applications.

### 3.4. NeRFing Corrective Noise Augmentation

Given a visual augmentation model (Sec. 3.3), we can adapt methods for corrective augmentation (Sec. 3.1) into the visual domain. Our goal is to create noise-distributed corrective transitions  $\{(\tilde{I}_k, \tilde{T}_k, \tilde{a}_k)\}$ . First, we describe this simply in the global frame. In order to sample from the noise-distributed corrective trajectory, one can first apply noise to the measured end-effector pose,  $\tilde{T}_k := T_k \varepsilon$ , where  $\varepsilon \sim \text{NoiseDist}(SE(3))$  is a randomly sampled rotation and translation. The high-fidelity perturbed image  $\tilde{I}_k$  corresponding to  $\tilde{T}_k$  can then be rendered using the trajectory-specific NeRF  $F_{\Theta}^{\tau}$ , without requiring access to the actual environment. For the actions, the simplest case is when they

are inputs to a global-frame stabilizing Cartesian controller [25], in which case  $\tilde{a}_k = a_k = \tilde{T}_k$  will provide stabilization to the nominal trajectory, where  $\tilde{T}_k$  is the *desired* pose sent to the lower-level controller.

**Corrective Relative Actions.** As is common in prior works, we observe better performance by parameterizing the learned policy as a *relative* rather than global action. Consider the as-discussed global-frame version, with (1) *observations* as a global-frame measured SE(3) end-effector pose  ${}^W T_k^E$ , where  $W$  refers to world-frame, and  $E$  to the end-effector frame at timestep  $k$ , and (2) *action* as a global-frame desired SE(3) end-effector pose  ${}^W T_k^{\hat{E}}$ . To switch to a relative action space, we adjust the action to

$$a_k = {}^E T_k^{\hat{E}} = ({}^W T_k^E)^{-1} {}^W T_k^{\hat{E}} = {}^E T_k^W {}^W T_k^{\hat{E}}. \quad (4)$$

To additionally formulate the corrective noise augmentation in the relative frame, we consider the SE(3)-noise  $\varepsilon$  as transforming from the noisy end-effector frame to the measured end-effector frame, i.e.  ${}^E T_k^{\hat{E}} := \varepsilon$ . This accordingly adjusts the observation as  ${}^W T_k^{\hat{E}} = {}^W T_k^E {}^E T_k^{\hat{E}} = {}^W T_k^E \varepsilon$  and the *relative* action as:

$$\tilde{a}_k = {}^{\hat{E}} T_k^{\hat{E}} = ({}^W T_k^E {}^E T_k^{\hat{E}})^{-1} {}^W T_k^{\hat{E}} = \varepsilon^{-1} a_k. \quad (5)$$

Concisely, this amounts to post-pending  $\varepsilon$  to the measured pose, and pre-pending  $\varepsilon^{-1}$  to the un-noised relative action. **Non-SE(3) actions.** Thus far we have assumed that the actions  $a \in SE(3)$  only command desired pose transforms. In practice, the action space may contain *additional* dimensions, for example to open and close the gripper itself. The SPARTN augmentation does not concern these aspects of the action, so the corrective actions simply copy the original action for non-SE(3) dimensions.

**Summary.** Figure 2 summarizes our augmentation procedure: given a demonstration dataset  $\mathcal{D}$ , we first train all the neural radiance fields  $\{F_{\Theta}^{\tau}\}_{\tau \in \mathcal{D}}$  and save the weights to disk. We then augment each transition in the original dataset with  $N_{\text{aug}}$  noisy corrective transitions produced using the process we described above, and save these transitions into an augmented dataset  $\tilde{\mathcal{D}}$ . Appendix Algorithm 2 describes the precise augmentation procedure in detail.

After the augmented dataset has been created, it can simply be combined with the original dataset to augment BC training. Various sampling strategies are possible, but in our experiments we simply construct mini-batches for BC training by sampling from the original and augmented datasets with equal probability.

## 4. Simulation Experiments

We evaluate SPARTN and related approaches in the simulated 6-DoF grasping benchmark first introduced in [58], which features a simulated Franka Emika Panda arm with

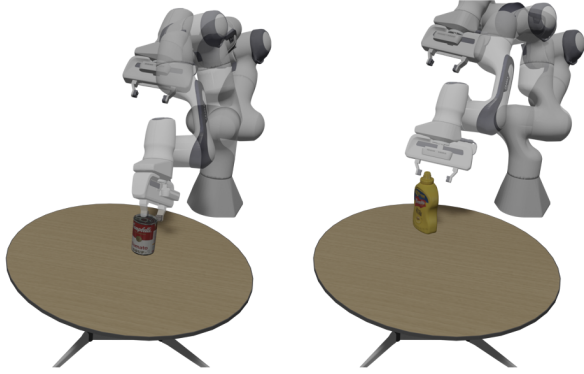


Figure 5. Example tasks in the simulated 6-DoF grasping benchmark, first introduced in [58]. There are  $\sim 1,500$  ShapeNet objects in training demonstrations, and held out YCB objects for evaluation. A camera mounted to the Franka Panda robot’s arm provides observations, and the policy controls the 6-DoF end-effector pose.

a parallel-jaw gripper. Objects are placed on a table and must be lifted above a height threshold in a successful grasp. Policies receive either RGB, RGBD, or point cloud observations from a wrist camera, and control the gripper by commanding relative pose changes in 6-DoF end-effector space.

#### 4.1. Data Collection and Evaluation Protocol

We follow the training and evaluation procedure from [58]. The training dataset includes 2,500 demonstrations of grasping 1,500 ShapeNet [6] objects. Demonstrations are up to 20 timesteps and are generated by trajectory optimization to precomputed grasps from the ACRONYM dataset [11]. Policies are evaluated on grasping *held-out* objects from the YCB [5] or ShapeNet datasets. Though it is more “out of distribution” relative to the training objects, the YCB evaluation is more realistic for tabletop grasping scenarios (ShapeNet includes a wide variety of objects, e.g. airplanes and bicycles). Each held-out object is evaluated ten times, each with a different initial robot configuration and object’s initial pose.

#### 4.2. Comparisons

We compare SPARTN against other approaches ranging from simple behavior cloning to online-supervised imitation and reinforcement learning:

**Behavior Cloning (BC):** Trains policies via supervised learning on the demonstration dataset.

**DART [29]:** Introduces a modified demonstration setup where a continuous Gaussian noise is injected into the robot’s state *during* the expert demonstration, then trains policies on the modified demonstrations with BC.

**Homography Augmentation (HA):** A simplification of SPARTN where perturbations can be 3D rotations, but not

<sup>1</sup>Note that the “BC” results in [58] actually use DART.

Supervision	Method	Input	YCB SR(%)	SN SR(%)
Offline	BC	RGB	28.9 $\pm$ 2.4	57.4 $\pm$ 0.2
		RGB	51.2 $\pm$ 3.2	57.8 $\pm$ 2.2
	DART <sup>†</sup>	RGBD	46.3	45.3
		Point cloud	65.6	73.6
	HA	RGB	29.7 $\pm$ 2.4	57.5 $\pm$ 0.8
	<b>SPARTN (ours)</b>	RGB	<b>74.7 <math>\pm</math> 2.4</b>	<b>66.9 <math>\pm</math> 1.6</b>
Online*		RGB	52.3	52.1
	Dagger	RGBD	67.1	60.4
		Point cloud	77.2	75.8
	GA-DDPG	Point cloud	88.2	91.3

Table 2. Grasping success rates (SR) on *held-out objects* from YCB [5] or ShapeNet (SN) [6] in a simulated 6-DoF grasping benchmark [58]. We **bold** the best offline RGB-only results, though we include online and non-RGB methods for comparison. Online\* requires additional environment interactions, while Offline only uses demonstration data. DART<sup>†</sup> is offline but requires a special demonstration collection setup. SPARTN outperforms other offline RGB-only methods, while RL (GA-DDPG) performs best overall while requiring millions of interactions. We calculate average success rates and standard error over 4 random seeds. *Italicized* success rates were reported in prior work<sup>1</sup> [58].

Method	Image aug.	YCB SR(%)
BC	Without	25.3 $\pm$ 1.6
	With	28.9 $\pm$ 2.4
DART	Without	51.2 $\pm$ 3.2
	With	48.6 $\pm$ 2.8
HA	Without	23.9 $\pm$ 4.0
	With	29.7 $\pm$ 2.4
SPARTN	Without	68.3 $\pm$ 3.6
	With	<b>74.7 <math>\pm</math> 2.4</b>

Table 3. Ablating the effect of standard image augmentation on each method for RGB policies. Average success rates and standard errors are calculated over four seeds.

translations, of the camera. For pure rotation, we can calculate the homography transform for rendering the rotated view without NeRF. Augmented actions are computed similarly to SPARTN.

**Dagger [46]:** An online-supervised method where a policy is first trained using BC on offline demos, and then the expert provides action labels for states from the policy’s roll-outs throughout further policy training.

**GA-DDPG [58]:** Jointly trains policies via BC and fine-tunes them with reinforcement learning (DDPG [51]).

Of the methods considered, Dagger and GA-DDPG require online environment interaction and supervision from an expert or reward function, respectively. The other methods, including SPARTN, train only on pre-collected demonstration datasets. Since the DART data collection setup noisily perturbs the robot as the expert demonstrates, DART can interfere with human expert demonstration, though this challenge does not arise here with our simulated expert.

### 4.3. Training Details

We follow the RGB architectural and training hyperparameters of [58], with policies consisting of a ResNet-18 [15] image encoder followed by an MLP. We apply random crop and color jitter to the training images. We re-use the same training BC hyperparameters for SPARTN. Appendix C.1 describes the complete training details.

For SPARTN, we create  $N_{aug} = 100$  augmented transitions from each original transition and save the augmented dataset to disk before BC training. To sample the perturbations  $\varepsilon \sim \text{NoiseDist}(\text{SE}(3))$ , we parameterize the rotations in terms of Euler angles  $(\phi, \theta, \varphi)$  and uniformly sample both rotation and translation parameters:

$$(\phi, \theta, \varphi), (t_x, t_y, t_z) \sim \mathcal{U}(-\alpha, \alpha), \mathcal{U}(-\beta, \beta) \quad (6)$$

$$\varepsilon := (R(\phi, \theta, \varphi), (t_x, t_y, t_z)) \quad (7)$$

In simulation, we set  $\alpha = 0.2$  radians and  $\beta = 3$  mm. Following the DART hyperparameters in [58], we only augment timesteps 5 – 13 of each demonstration. Appendix B contains samples of SPARTN’s perturbed observations rendered via NeRF.

### 4.4. Results

Table 2 shows grasping success rates on held-out objects from either the YCB or ShapeNet (SN) datasets. SPARTN significantly outperforms the other two offline-supervised approaches for RGB inputs. Since DART injects noise during expert collection, the amount of augmentation is limited by the number of demonstrations the expert can collect. Meanwhile, SPARTN can cheaply generate an arbitrary number of augmented examples on top of the existing demonstration dataset, leading to more data diversity without any additional effort from the expert. See the supplementary material for videos of rollouts for SPARTN and BC policies, with success and failure cases.

Naive BC performs relatively well on ShapeNet evaluation, likely because the evaluation set is “in-distribution”. Correspondingly, corrective methods (DART, SPARTN) achieve smaller improvements over BC on ShapeNet and larger improvements on YCB.

On YCB, SPARTN’s performance is closer to RL (GADDPG) than to other offline methods. It actually outperforms DAgger for RGB policies, perhaps because DAgger’s data distribution changes online presenting a more challenging learning problem. Notably, on YCB SPARTN outperforms DART with point clouds and is comparable to DAgger with point clouds. Since SPARTN itself only requires RGB images during both NeRF training and policy training, it significantly closes the gap between the performance of RGB-only and the best depth-based methods. Since consumer depth cameras can struggle with common thin and reflective items [62], improving RGB-only grasping can enable more robust grasping of such objects.

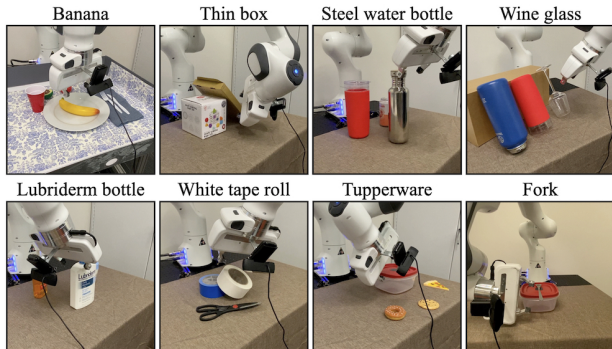


Figure 6. Real-world grasping environments. In each environment, the task is to grasp the labeled target object in a particular way. The target objects have various geometric shapes and exhibit a diverse range of characteristics, including reflectiveness, transparency, and radial symmetry.

Target Object	# Demos	BC SR(%)	SPARTN SR(%)
Banana	14	55	<b>75</b>
Thin box	20	35	<b>65</b>
Steel water bottle	15	20	<b>40</b>
Wine glass	25	70	<b>90</b>
Lubriderm bottle	17	25	<b>60</b>
White tape roll	15	40	<b>45</b>
Tupperware	20	40	<b>75</b>
Fork	20	25	<b>40</b>
Average	--	38.75	<b>61.25</b>

Table 4. Success rates (SR) of behavior cloning (BC) and SPARTN 6-DoF grasping policies on a suite of eight real-world target objects. SPARTN outperforms BC in every environment, achieving an average absolute performance boost of 22.5%. Each success rate is computed over 20 trials.

**Effect of Image Augmentations.** We try ablating the effect of standard image augmentations (random crop and color jitter) for the BC, DART, HA, and SPARTN methods. Table 3 shows that the image augmentations have a small effect on the performance of most methods, relative to the larger effect of SPARTN’s corrective augmentations.

## 5. Real-World Experiments

The simulation benchmark results show that SPARTN can improve grasping generalization in imitation learning without online supervision. Here, we verify that SPARTN can enable real-world robotic grasping of challenging objects from limited human demonstration. See the website in the supplement for video results.

### 5.1. Experimental details

**Robot setup.** The robotic manipulator is a Franka Emika Panda robot arm with a wrist-mounted consumer grade webcam. Policies take images of size  $256 \times 256$  as input and output a 6-DoF action representing the desired change in end-effector position and orientation, as well as a binary

open/close gripper action. We use a Cartesian impedance controller to command the pose changes at a frequency of 4 Hz. We task the robot to grasp a target object in eight different environments depicted in Figure 6. The target objects include natural shapes that are common in the real world and exhibit a diverse range of attributes, such as reflectiveness, transparency, and radial symmetry.

**Comparisons and Evaluation.** In each environment, we collect a small number of expert grasping demonstrations with a virtual reality controller. Because DART is difficult to use with human demonstration collection, we compare SPARTN to a vanilla BC policy trained on the same set of demonstrations. Policies are evaluated with the same objects seen in the demonstrations. Initial object and robot configurations are randomized during both data collection and evaluation.

**Training.** To improve NeRF quality for SPARTN, we program the robot to collect a few images of the scene from a fixed set of poses before the start of each demonstration. This automatic process is only used to improve COLMAP’s pose estimation and the subsequent NeRF training. For SPARTN, we generate  $N_{aug} = 50$  augmented transitions from each original transition in the demonstrations. We sample perturbations  $\epsilon$  according to Eq. 6 with  $\alpha = 0.05$  radians and  $\beta = 0.4$  mm. Appendix D.2 describes COLMAP and NeRF training in more detail, and Appendix B shows sample SPARTN observations rendered by NeRF. Aside from using the augmented dataset, SPARTN policies are trained using the same BC architecture and training hyperparameters described in Appendix D.1.

## 5.2. Results

Table 4 shows grasping success rates in the eight real-world environments. Quantitatively, SPARTN policies outperform the baseline BC policies across the board, on average achieving an absolute 22.5% increase in success rate.

Figure 7 shows qualitative differences in performance between the BC and SPARTN policies. SPARTN generally exhibits more reactive behaviors than the baseline policy: it navigates towards the target object better while occasionally avoiding obstacles, executes actions with greater precision, and even reattempts the grasp more successfully after an initial miss. In some cases, the differences are stark: for instance, SPARTN may successfully move toward and grasp the target object while the baseline fails to even reach it. We present further analysis of policy rollouts in Appendix D.3, showing how SPARTN qualitatively performs better than BC even in cases where both methods fail to grasp the target object. The supplementary videos of real-world policy rollouts illustrate all of these differences, revealing that SPARTN induces important reactive closed-loop behaviors that enable the manipulator to successfully execute grasps in the real world.

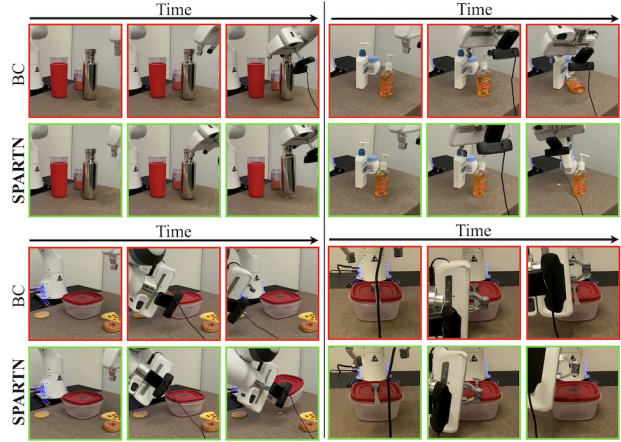


Figure 7. Sample real-world policy rollouts illustrating how SPARTN succeeds (green) in cases where BC fails (red). **Top left:** While BC fails to reach the steel bottle, SPARTN successfully reaches and grasps it. **Top right:** While BC collides into the orange bottle (a distractor object), SPARTN takes a more rounded path to avoid it before grasping the white Lubriderm bottle. **Bottom left:** BC fails to recover after a missed grasp, while SPARTN successfully reattempts the grasp after failing the first time. **Bottom right:** SPARTN operates with higher precision than BC and successfully completes a difficult fork grasping task.

## 6. Conclusion

We introduce SPARTN, which augments eye-in-hand demonstrations with perturbed visual observations and corrective actions. Our augmentation leverages novel-view synthesis, in particular NeRF, to produce these augmentations during *training*-time. SPARTN can improve behavior cloning training of robust, real-time, and closed-loop 6-DoF visual control policies. We show that SPARTN-trained policies outperform other offline-supervised methods in a simulated 6-DoF grasping generalization benchmark. Our policies can also perform on par with imitation methods that require depth information and online supervision. We verify that SPARTN can train policies to grasp a variety of objects in the real world from limited human demonstrations.

Despite its strong performance, SPARTN also has some limitations. First, SPARTN is limited to tasks with static scenes like grasping: extending to a broader set of manipulation tasks would require effective view synthesis for dynamic scenes. Second, training a neural radiance field for every demonstration before training is computationally expensive, a limitation that may be mitigated through amortized NeRF models [55, 59, 64]. Finally, the noise distribution must be tuned for a given platform, e.g. it is tuned separately for the simulated and real experiments. An interesting direction for future work is to use policy actions to generate the noise, which would result in a fully offline variant of DAGger [46] that uses NeRF as a simulator.



## 7. Acknowledgements

We thank Jimmy Wu, Kaylee Burns, Bohan Wu, and Suraj Nair for technical advice on various aspects of the real robot setup, and Archit Sharma for helpful conceptual discussions. This project was funded by ONR grant N00014-21-1-2685. AZ acknowledges the support of the NSF Graduate Research Fellowship. CF is a fellow of the CIFAR Learning in Machines and Brains Program.

## References

- [1] Michal Adamkiewicz, Timothy Chen, Adam Caccavale, Rachel Gardner, Preston Culbertson, Jeannette Bohg, and Mac Schwager. Vision-only robot navigation in a neural radiance world. *IEEE Robotics and Automation Letters*, 7(2):4606–4613, 2022. 3
- [2] Antonio Bicchi and Vijay Kumar. Robotic grasping and contact: A review. In *Proceedings 2000 ICRA. Millennium conference. IEEE international conference on robotics and automation. Symposia proceedings (Cat. No. 00CH37065)*, volume 1, pages 348–353. IEEE, 2000. 2
- [3] Valts Blukis, Taeyeop Lee, Jonathan Tremblay, Bowen Wen, In So Kweon, Kuk-Jin Yoon, Dieter Fox, and Stan Birchfield. Neural fields for robotic object manipulation from a single image. *arXiv preprint arXiv:2210.12126*, 2022. 3
- [4] Jeannette Bohg, Antonio Morales, Tamim Asfour, and Danica Kragic. Data-driven grasp synthesis—a survey. *IEEE Transactions on robotics*, 30(2):289–309, 2013. 2
- [5] Berk Calli, Aaron Walsman, Arjun Singh, Siddhartha Srinivasa, Pieter Abbeel, and Aaron Dollar. Benchmarking in manipulation research: The YCB object and model set and benchmarking protocols. *IEEE Robotics and Automation Magazine*, 22(3):36–52, 2015. 6
- [6] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An information-rich 3d model repository. *arXiv:1512.03012*, 2015. 6
- [7] Simon Le Cleac’h, Hong-Xing Yu, Michelle Guo, Taylor A Howell, Ruohan Gao, Jiajun Wu, Zachary Manchester, and Mac Schwager. Differentiable physics simulation of dynamics-augmented neural objects. *arXiv preprint arXiv:2210.09420*, 2022. 3
- [8] Mehmet Dogar and Siddhartha Srinivasa. A framework for push-grasping in clutter. *Robotics: Science and systems VII*, 1, 2011. 1
- [9] Danny Driess, Zhiao Huang, Yunzhu Li, Russ Tedrake, and Marc Toussaint. Learning multi-object dynamics with compositional neural radiance fields. *arXiv preprint arXiv:2202.11855*, 2022. 3
- [10] Danny Driess, Ingmar Schuber, Pete Florence, Yunzhu Li, and Marc Toussaint. Reinforcement learning with neural radiance fields. *arXiv preprint arXiv:2206.01634*, 2022. 3
- [11] Clemens Eppner, Arsalan Mousavian, and Dieter Fox. ACRONYM: A large-scale grasp dataset based on simulation. *arXiv:2011.09584*, 2020. 6
- [12] Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: A large-scale benchmark for general object grasping. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11444–11453, 2020. 2
- [13] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. In *International conference on machine learning*, pages 49–58. PMLR, 2016. 1
- [14] Peter Florence, Lucas Manuelli, and Russ Tedrake. Self-supervised correspondence in visuomotor policy learning. *IEEE Robotics and Automation Letters*, 5(2):492–499, 2019. 1, 2, 3
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *corr abs/1512.03385* (2015), 2015. 7, 12
- [16] Ryan Hoque, Ashwin Balakrishna, Ellen Novoseller, Albert Wilcox, Daniel S Brown, and Ken Goldberg. Thriftydagger: Budget-aware novelty and risk gating for interactive imitation learning. In *5th Annual Conference on Robot Learning*, 2021. 2
- [17] Kyle Hsu, Moo Jin Kim, Rafael Rafailov, Jiajun Wu, and Chelsea Finn. Vision-based manipulators need to also see from their hands. *arXiv preprint arXiv:2203.12677*, 2022. 2
- [18] Jeffrey Ichnowski, Yahav Avigal, Justin Kerr, and Ken Goldberg. Dex-nerf: Using a neural radiance field to grasp transparent objects. *arXiv preprint arXiv:2110.14217*, 2021. 2, 3
- [19] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning*, pages 991–1002. PMLR, 2022. 1
- [20] Rishabh Jangir, Nicklas Hansen, Sambaran Ghosal, Mohit Jain, and Xiaolong Wang. Look closer: Bridging egocentric and third-person views with transformers for robotic manipulation. *IEEE Robotics and Automation Letters*, 7(2):3046–3053, 2022. 2
- [21] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pages 651–673. PMLR, 2018. 1, 2
- [22] Liyiming Ke, Jingqiang Wang, Tapomayukh Bhattacharjee, Byron Boots, and Siddhartha Srinivasa. Grasping with chopsticks: Combating covariate shift in model-free imitation learning for fine manipulation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6185–6191. IEEE, 2021. 1, 2, 3
- [23] Michael Kelly, Chelsea Sidrane, Katherine Driggs-Campbell, and Mykel J Kochenderfer. Hg-dagger: Interactive imitation learning with human experts. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8077–8083. IEEE, 2019. 2
- [24] Justin Kerr, Letian Fu, Huang Huang, Jeffrey Ichnowski, Matthew Tancik, Yahav Avigal, Angjoo Kanazawa, and Ken

- Goldberg. Evo-nerf: Evolving nerf for sequential robot grasping. In *6th Annual Conference on Robot Learning*, 2022. 2, 3
- [25] Oussama Khatib. A unified approach for motion and force control of robot manipulators: The operational space formulation. *IEEE Journal on Robotics and Automation*, 3(1):43–53, 1987. 5
- [26] Kilian Kleeberger, Richard Bormann, Werner Kraus, and Marco F Huber. A survey on learning-based robotic grasping. *Current Robotics Reports*, 1(4):239–249, 2020. 2
- [27] Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *arXiv preprint arXiv:2004.13649*, 2020. 3
- [28] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054*, 2022. 12
- [29] Michael Laskey, Jonathan Lee, Roy Fox, Anca Dragan, and Ken Goldberg. Dart: Noise injection for robust imitation learning. In *Conference on robot learning*, pages 143–156. PMLR, 2017. 2, 6
- [30] Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. *Advances in neural information processing systems*, 33:19884–19895, 2020. 2, 3
- [31] Ian Lenz, Honglak Lee, and Ashutosh Saxena. Deep learning for detecting robotic grasps. *The International Journal of Robotics Research*, 34(4-5):705–724, 2015. 2
- [32] Sergey Levine, Peter Pastor, Alex Krizhevsky, Julian Ibarz, and Deirdre Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International journal of robotics research*, 37(4-5):421–436, 2018. 2
- [33] Yunzhu Li, Shuang Li, Vincent Sitzmann, Pulkit Agrawal, and Antonio Torralba. 3d neural scene representations for visuomotor control. In *Conference on Robot Learning*, pages 112–123. PMLR, 2022. 3
- [34] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio Ojea, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. *arXiv preprint arXiv:1703.09312*, 2017. 1, 2
- [35] Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. *Conference on Robot Learning (CoRL)*, 2021. 2
- [36] Kunal Menda, Katherine Driggs-Campbell, and Mykel J Kochenderfer. Ensembledagger: A bayesian approach to safe imitation learning. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5041–5048. IEEE, 2019. 2
- [37] 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. *Communications of the ACM*, 65(1):99–106, 2021. 2, 3
- [38] Andrew T Miller and Peter K Allen. Graspit! a versatile simulator for robotic grasping. *IEEE Robotics & Automation Magazine*, 11(4):110–122, 2004. 1
- [39] Douglas Morrison, Peter Corke, and Jürgen Leitner. Closing the loop for robotic grasping: A real-time, generative grasp synthesis approach. In *Robotics: Science and Systems (RSS)*, 2018. 2
- [40] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof graspnet: Variational grasp generation for object manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2901–2910, 2019. 1, 2
- [41] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 3, 12
- [42] Adithyavairavan Murali, Arsalan Mousavian, Clemens Eppner, Chris Paxton, and Dieter Fox. 6-dof grasping for target-driven object manipulation in clutter. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6232–6238. IEEE, 2020. 2
- [43] Lerrel Pinto and Abhinav Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In *2016 IEEE international conference on robotics and automation (ICRA)*, pages 3406–3413. IEEE, 2016. 2
- [44] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv:1709.10087*, 2017. 1
- [45] Joseph Redmon and Anelia Angelova. Real-time grasp detection using convolutional neural networks. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 1316–1322, 2015. 2
- [46] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2011. 1, 2, 6, 8
- [47] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 12
- [48] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4, 12
- [49] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 4, 12
- [50] Karun B Shimoga. Robot grasp synthesis algorithms: A survey. *The International Journal of Robotics Research*, 15(3):230–266, 1996. 2
- [51] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy

- gradient algorithms. In *International conference on machine learning*, pages 387–395. PMLR, 2014. 6
- [52] Anthony Simeonov, Yilun Du, Andrea Tagliasacchi, Joshua B Tenenbaum, Alberto Rodriguez, Pulkit Agrawal, and Vincent Sitzmann. Neural descriptor fields: Se (3)-equivariant object representations for manipulation. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 6394–6400. IEEE, 2022. 2
- [53] Shuran Song, Andy Zeng, Johnny Lee, and Thomas Funkhouser. Grasping in the wild: Learning 6dof closed-loop grasping from low-cost demonstrations. *IEEE Robotics and Automation Letters*, 5(3):4978–4985, 2020. 2
- [54] Andreas ten Pas, Marcus Gualtieri, Kate Saenko, and Robert Platt. Grasp pose detection in point clouds. *The International Journal of Robotics Research*, 36(13-14):1455–1473, 2017. 1, 2
- [55] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d representation and rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15182–15192, 2021. 8
- [56] Lirui Wang, Xiangyun Meng, Yu Xiang, and Dieter Fox. Hierarchical policies for cluttered-scene grasping with latent plans. *IEEE Robotics and Automation Letters*, 7(2):2883–2890, 2022. 2
- [57] Lirui Wang, Yu Xiang, and Dieter Fox. Manipulation trajectory optimization with online grasp synthesis and selection. *arXiv preprint arXiv:1911.10280*, 2019. 1, 2
- [58] Lirui Wang, Yu Xiang, Wei Yang, Arsalan Mousavian, and Dieter Fox. Goal-auxiliary actor-critic for 6d robotic grasping with point clouds. In *Conference on Robot Learning*, pages 70–80. PMLR, 2022. 1, 2, 5, 6, 7, 12
- [59] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 8
- [60] Xinchun Yan, Jasmined Hsu, Mohammad Khansari, Yunfei Bai, Arkanath Pathak, Abhinav Gupta, James Davidson, and Honglak Lee. Learning 6-dof grasping interaction via deep geometry-aware 3d representations. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3766–3773. IEEE, 2018. 1
- [61] Denis Yarats, Ilya Kostrikov, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In *International Conference on Learning Representations*, 2020. 2
- [62] Lin Yen-Chen, Pete Florence, Jonathan T Barron, Tsung-Yi Lin, Alberto Rodriguez, and Phillip Isola. Nerf-supervision: Learning dense object descriptors from neural radiance fields. *arXiv preprint arXiv:2203.01913*, 2022. 3, 7
- [63] Lin Yen-Chen, Pete Florence, Jonathan T Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. inerf: Inverting neural radiance fields for pose estimation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1323–1330. IEEE, 2021. 3
- [64] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 8

## A. Full Algorithm

---

### Algorithm 2 SPARTN: Corrective Augmentation via NeRF

---

**Input:**  $\mathcal{D} = \{\tau\}_{i=1}^N$  - expert demonstrations  
**Input:**  $E T^C$  - Transform between camera and end-effector  
**Input:**  $M \in \{0, 1\}^{w \times h}$  - mask for robot gripper  
**Output:** Augmented transition dataset  $\tilde{\mathcal{D}}$

- 1:  $\tilde{\mathcal{D}} \leftarrow \{\}$
- 2: **for**  $\tau = \{(I_k, {}^W T_k^E), ({}^E T_k^{\tilde{E}})\}_{k=1}^K \in \mathcal{D}$  **do**
- 3:   **# EE pose  $\rightarrow$  camera pose**
- 4:    $W T_k^C \leftarrow W T_k^E E T^C$ , for  $k = 0, \dots, K$
- 5:    $\{V T_k^C\}_{k=1}^K \leftarrow \text{COLMAP}(I_0, \dots, I_K)$
- 6:    $\beta \leftarrow \text{SOLVE}(\text{Eq. 8})$  **# Estimate  $\beta$**
- 7:    $V T_k^W = V T_k^C ({}^W T_k^C)^{-1}$ , for  $k = 0, \dots, K$
- 8:   **# train NeRF**
- 9:    $F_{\Theta}^{\tau} \leftarrow \text{NeRF}(\{(I_k, V T_k^C)\}_{k=1}^K, M)$
- 10:   **for**  $k \in (0, 1, \dots, K)$  **do**
- 11:     **for**  $i = 1 : N_{\text{aug}}$  **do**
- 12:        $\varepsilon \leftarrow \text{NoiseDist}(\text{SE}(3))$
- 13:        $W T_k^{\tilde{E}} \leftarrow W T_k^E \varepsilon$  **# perturbed end-effector pose**
- 14:        $\tilde{E} T_k^{\tilde{E}} \leftarrow \varepsilon^{-1} E T_k^{\tilde{E}}$  **# corrective relative action**
- 15:        $W T_k^{\tilde{C}} \leftarrow W T_k^{\tilde{E}} E T^C$  **# perturbed camera pose**
- 16:        $V T_k^{\tilde{C}} \leftarrow V T_k^W W T_k^{\tilde{C}}$
- 17:        $\tilde{I}_k \leftarrow \text{NeRF-Render}(V T_k^{\tilde{C}}; F_{\Theta}^{\tau})$
- 18:       **# splice gripper into rendered image**
- 19:        $\tilde{I}_k \leftarrow \neg M \odot \tilde{I}_k + M \odot I_k$
- 20:        $\tilde{\mathcal{D}} \leftarrow \tilde{\mathcal{D}} \cup \{(\tilde{I}_k, {}^W T_k^{\tilde{E}}), ({}^{\tilde{E}} T_k^{\tilde{E}})\}$

---

## B. Sample Perturbed Observations

As described in the main text, image observations  $\tilde{I}$  for perturbed camera poses are rendered via novel-view synthesis techniques, in particular NeRF. The section ‘‘Example NeRF renders from augmented poses’’ in the supplementary website shows videos of the rendering output for both simulation and real-world augmentation.

## C. Simulation Details

### C.1. Training for Simulation Experiments

We train all of our RGB policies (SPARTN, DART, and HA) closely following the architectural and optimization hyperparameters of [58]. The image encoder is a ResNet-18 [15] pretrained on ImageNet [47], followed by a 3-layer MLP with 512 hidden units each and ReLU activations. The MLP outputs the predicted action as a 6-D vector of predicted translations and Euler angles. The final layer has a tanh activation to scale the output between  $[-1, 1]$ , then each dimension is scaled to match the environment’s action bounds.

For behavior cloning, we use the same 3D point-matching loss as [58]. We optimize the objective using the Adam optimizer with batch size 100 for 100,000 steps. We train only the MLP (and not the pre-trained ResNet) for the first 20,000 steps, since ‘‘freeze-then-train’’ fine-tuning techniques have been shown to be more robust [28]. After 20,000 steps, the ResNet is unfrozen and the entire network is trained as usual.

### C.2. NeRF Training

We use Instant-NGP [41] to train a NeRF for each of the 2,500 demonstration scenes. We train each NeRF for 3,500 steps, which takes 30 seconds on an NVIDIA GeForce RTX 2080 Ti. To reduce total training time, we train the NeRFs in parallel: with 4 GPUs, we can train all 2,500 NeRFs in  $\sim 7$  hours. Because camera calibration is exact in simulation, we use the world-frame camera poses given by calibration instead of COLMAP to train NeRF.

## D. Real-World Details

### D.1. Policy Training

For both BC and SPARTN, the image encoder for each policy is a four-layer convolutional network followed by two feedforward layers, and uses batch normalization and ReLU activation functions. The image embedding is fed into a three-layer ReLU MLP policy head, which outputs the predicted action. We train the policy from scratch with mean squared error, using Adam with learning rate  $5 \times 10^{-4}$  and a batch size of 64. No image augmentations like random crop are applied.

### D.2. COLMAP and NeRF Training

In the real world, we train our NeRF using camera poses  $\{V H_k^C\}_{k=1}^K$  estimated by COLMAP [48, 49]. From the camera calibration, we also have pose transforms  $\{W T_k^C\}_{k=1}^K$  which are too noisy to train a good NeRF from, but will allow us to convert between world and COLMAP frames. Conversion is necessary because we want to render from perturbed camera poses  $W T_k^{\tilde{C}}$ , but we must call NeRF-Render using the transform  $V H_k^{\tilde{C}}$  (COLMAP’s reference frame) instead. From Eq. 2, we can do this conversion if we have the transform  $V T_k^W$  and the scale factor  $\beta$  which distinguishes  $T$  and  $H$ .

We can estimate both quantities using the  $(W T_k^C, V H_k^C)$  pairs in each demonstration. Recalling that  ${}^a T^b = ({}^b T^a)^{-1}$ , we have for any  $j \neq k$ :

$$\beta \text{Trans} [{}^C T_j^W W T_k^C] = \text{Trans} [{}^C H_j^V V H_k^C] \quad (8)$$

This leads to an overconstrained linear system in one unknown  $\beta$ , which can be solved with linear regression. After solving for  $\beta$ , we can rescale from  $H \rightarrow T$ :

$$V T_k^C := (\text{Rot} [{}^C H_k^V], (1/\beta) \text{Trans} [{}^C H_k^V]) \quad (9)$$

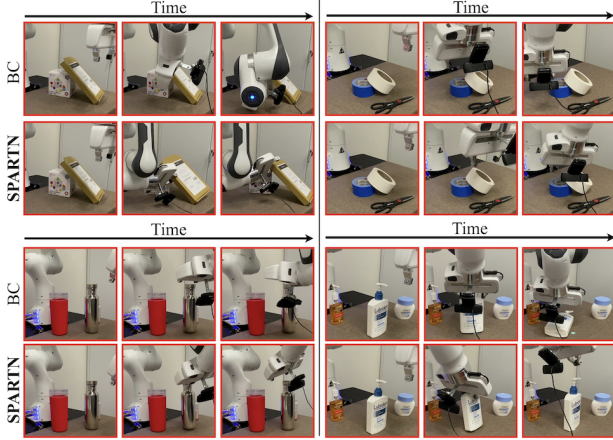


Figure 8. Sample real-world policy rollouts illustrating cases where both BC and SPARTN policies fail. **Top left:** BC fails to reach the thin box, while SPARTN successfully reaches it and nearly completes the grasp. **Top right:** While BC completely misses the white tape roll, SPARTN makes contact but fails to grasp it. **Bottom left:** BC fails to reach the steel water bottle, while SPARTN contacts it but fails to grasp it. **Bottom right:** BC reaches the Lubriderm bottle over but knocks it over before grasping it, while SPARTN nearly completes the grasp as the bottle barely slips away from the robot’s fingers.

Now solving for  ${}^V T_k^W$  is simple:

$${}^V T_k^W = {}^V T_k^C ({}^W T_k^C)^{-1}. \quad (10)$$

In theory,  ${}^V T_k^W$  is constant for all  $k$ , but in practice there is noise from both COLMAP and the real-world pose estimates, so we simply compute  ${}^V T_k^W$  separately for each  $k$ .

### D.3. Additional Qualitative Analysis

Figure 8 shows sample evaluation trials where both SPARTN and BC policies fail, given the same initial object and end-effector configurations. Even in these failure cases, SPARTN qualitatively performs better than BC, nearly grasping the object in some cases, and dropping the object soon after grasping in other cases. See the videos on the supplemental website for additional qualitative results.

### D.4. Real-World Environments

The real-world environments shown in Figure 6 each include a single target object to grasp among other distractor objects. The initial positions of all objects and the robotic end-effector are randomized while collecting expert demonstrations and during test time. The initial orientations of the objects are also randomized such that all policies must perform end-effector orientation control to complete the task. To make a fair comparison, we evaluate different policies given the same set of initial configurations (as shown in Figure 7, Figure 8, and the videos on the supplemental website); we do this by resetting the robot and environment to

the same conditions and alternating the rollouts of different policies before moving on to the next test configuration. Rollouts are considered successful if the robot grasps the target object and lifts it steadily into the air. We describe the details of each environment below:

- **Banana:** The environment includes an artificial banana resting on a white ceramic plate, a red cup, an artificial avocado, an artificial tomato, and a plastic fork and knife resting on a blue cloth.
- **Thin box:** The environment includes a thin brown box resting on a thicker white toy box.
- **Steel water bottle:** The environment includes a steel water bottle, a red water bottle, and an aluminum can. The steel water bottle is highly reflective.
- **Wine glass:** The environment includes a wine glass, red water bottle, and blue water bottle all placed upside-down and leaning against a brown box. The wine glass is transparent.
- **Lubriderm bottle:** The environment includes a white Lubriderm bottle, an orange Neutrogena bottle, and a small white container of moisturizer.
- **White tape roll:** The environment includes a roll of white tape resting on a roll of blue tape, as well as a pair of black scissors. The rolls of tape are radially symmetric.
- **Tupperware:** The environment includes a plastic tupperware container with a red lid, as well as three artificial food items: a donut, a cookie, and a slice of pizza.
- **Fork:** The environment includes a fork resting on top of a plastic tupperware container. The fork’s location and orientation relative to the surface of the tupperware container is fixed, but the position of the entire assembly is randomized.