# HandNeRF: Learning to Reconstruct Hand-Object Interaction Scene from a Single RGB Image

Hongsuk Choi, Nikhil Chavan-Dafle, Jiacheng Yuan, Volkan Isler, and Hyunsoo Park
Samsung Research America, New York

*Abstract*— This paper presents a new method to learn hand-object interaction prior for reconstructing a 3D hand-object scene from a single RGB image. The inference as well as training-data generation for 3D hand-object scene reconstruction is challenging due to the depth ambiguity of a single image and occlusions by the hand and object. We turn this challenge into an opportunity by utilizing the hand shape to constrain the possible relative configuration of the hand and object geometry. We design a generalizable implicit function, HandNeRF, that explicitly encodes the correlation of the 3D hand shape features and 2D object features to predict the hand and object scene geometry. With experiments on real-world datasets, we show that HandNeRF can reconstruct hand-object scenes of novel grasp configurations more accurately than comparable methods. Moreover, we demonstrate that object reconstruction from HandNeRF ensures more accurate execution of downstream tasks, such as grasping and motion planning for robotic hand-over and manipulation. The code is released here: https://github.com/hongsukchoi/HandNeRF_RELEASE

## I. INTRODUCTION

The understanding of 3D hand-object interactions, i.e., semantic reconstruction of hand and object geometry, is key to applications such as human-to-robot object handover, and augmented and virtual reality. Most of the current methods are primarily based on template-based approaches, where a known 3D CAD model of a hand and an object is assumed. The major focus is predicting the transformation that fits the known 3D CAD model to input observation [1]–[3]. Even with these assumptions, the hand and object reconstruction from a single RGB image are both difficult tasks due to depth ambiguity, partial observation, and mutual occlusions.

The 3D hand reconstruction methods have seen significant advances [4]–[6] due to large-scale hand datasets and automated reliable 3D hand annotations [7]–[11]. In contrast, the progress in reconstruction of grasped objects from a single RGB image is relatively limited due to lack of data. Generating a 3D CAD model of large set of object and labeling 6D poses in hand-object interaction scenes are labor-intensive and challenging. The sparsity of views in realworld data collection settings makes the labeling ambiguous, often requiring manual initialization and post-processing for refining 6D object pose annotations [9], [12].

In this paper, we present a new method, named HandNeRF, that estimates a semantic neural radiance field of a hand-object interaction scene from a single RGB image and without using an object template. HandNeRF predicts the density (occupancy), color, and semantic label (hand, object, or background) for points in the 3D space which can be used for 3D semantic reconstruction and novel view synthesis. The key technical contribution of HandNeRF is that it alleviates
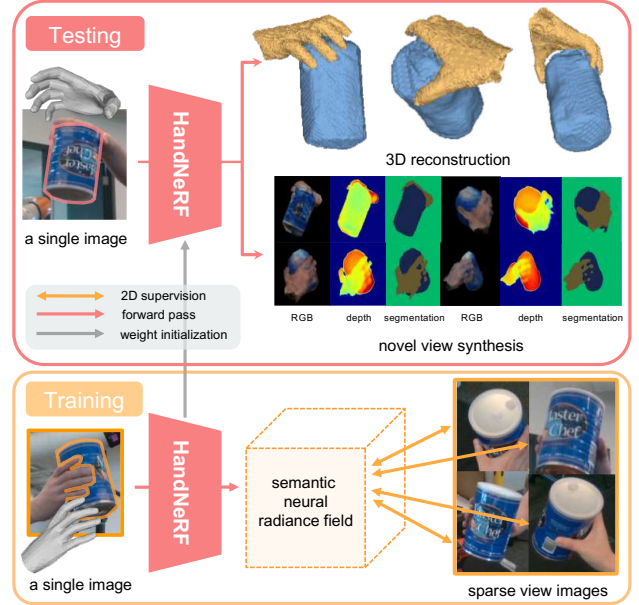


Fig. 1: Given a single RGB image of a hand-object interaction scene, HandNeRF predicts the hand and object's density, color, and semantics, which can be converted to reconstruction of 3D hand and object meshes and rendered to novel view images (RGB, depth, and semantic segmentation). HandNeRF learns the correlation between hand and object geometry from different types of hand-object interactions, supervised by sparse view images. HandNeRF is tested on a novel scene with an unseen hand-object interaction.

the ill-posed 2D to 3D reconstruction problem by utilizing the hand shape to constrain the possible relative configuration of hand and object geometry. In particular, HandNeRF explicitly learns the correlation between hand and object geometry to regularize their semantic reconstruction.

HandNeRF is trained on multiple hand-object interaction scenes to learn the correlation between hand and object geometry. Each scene has synchronized sparse view RGB images, 3D hand mesh annotation, and 2D semantic segmentation. At the inference time, a single RGB image with a novel grasp configuration is given. Fig. 1 illustrates HandNeRF, which is trained with sparse view RGB images and generates high-quality 3D reconstruction and rendering of images from an unseen single RGB input. Note that we do not use depth information in the whole process, which is much more unconstrained setting for both training and testing.

We evaluate HandNeRF on realworld datasets including

DexYCB [9] and HO-3D v3 [13] in terms of novel view synthesis and 3D reconstruction. We compare with the state-of-the-art template-free baselines [14]–[16], which we adapt to the task of reconstructing hand-object interaction without 3D object ground truth during training. Following the previous works [15], [17], we first keep the object used in training and testing the same, but the grasp configuration at testing is chosen to be significantly different from those during training. We further evaluate the generalization capability of HandNeRF by testing the model trained on 15 DexYCB objects on 4 unseen DexYCB objects. By learning the hand-object interaction prior with the explicit hand and object correlation, HandNeRF outperforms the baselines in generalization to novel hand grasps, which entail unseen occlusions and unseen object poses, and novel object shapes. Furthermore, we present qualitative results demonstrating HandNeRF's ability to reconstruct objects using in-house data. The annotation process for this data is fully automated in a casual environment, which uses only 7 sparse view RGB cameras, without the need for 3D CAD model generation or 6D object pose labeling. Finally, we experimentally demonstrate that HandNeRF enables more accurate execution of a downstream task of grasping for robotic hand-over and collision-free motion planning.

## II. RELATED WORK

Our work, HandNeRF, lies at the intersection of understanding 3D hand-object interaction and implicit neural representations. In this section, we first review the current approaches for 3D hand-object interaction reconstruction from a single RGB camera. Then, we discuss recent methods for sparse view-specific implicit neural representations, relevant to our work.

**3D reconstruction of hand-object interaction:** The study on the understanding of 3D hand-object interactions [12], [18]–[20] refers to semantic reconstruction of the hand and object geometry. In the context of this task, the existing methods for hand and object reconstruction are primarily based on template-based approaches, where the template indicates a known 3D CAD model of a hand and an object. The 3D hand reconstruction focuses on predicting mesh parameters, such as MANO [21], and has seen a significant advance due to large-scale datasets [7], [9], [10] and success of deep learning-based approaches [6], [22], [23]. Whereas, the 3D object reconstruction is approached as 6D pose estimation [1]–[3], which predicts the transformation that fits the known 3D CAD model to input observation.

The template-based approach for object reconstruction has two main limitations regarding collection of training data in the real world. First, it is costly and labor-intensive to obtain every object's 3D CAD model, requiring 3D laser scanning or a dense multi-view camera setup. Second, labeling 6D object poses in hand-object scenes is challenging due to hand occlusions and becomes more ambiguous if the captured views are not dense enough. In contrast, for training HandNeRF we require only a few sparse RGB views of hand-object interaction scenes and hand-pose annotations which

can be automated [10], [11].

Recently, [15], [19], [24] proposed methods that reconstruct a hand-held object without a known template. The work of Hasson et al. [19] jointly estimated the MANO hand parameters and genus-0 object mesh by leveraging AtlasNet [25]. Karunratanakul et al. [24] characterized the surface of the hand and object with a signed distance field. Ye et al. [15] conditioned object reconstruction on the hand articulation and also predicted a signed distance field of an object. While these methods are template-free at the inference time, they still require 3D object meshes for training. Therefore, they suffer with the same data collection problems as the template-based methods.

**Implicit neural representation from sparse view RGB images:** The sparse view-specific NeRF (Neural Radiance Field) representations have gained attention in object reconstruction [14], [26]–[28] and 3D human body reconstruction [16], [29]–[31]. Without any 3D annotation, they reconstruct a plausible 3D scene when optimized over a single scene only with sparse views. These methods address the limitations of multi-view reconstruction approaches such as vanilla NeRF [32] and SfM (Structure from Motion), which require a dense capture setup and fail when given sparse views [33]. These representations are explored for generalization by learning object appearance and geometry priors from multiple scenes. When tested on novel scenes with unseen object poses or unseen objects, a partial 3D reconstruction is achieved, although with blurry textures and noisy geometry. This limited performance is inevitable due to sparsity of input views, but the practical applications of these methods is significant. Nevertheless, scenes with a single view or hand-held objects are are less studied.

Our work is most relevant to the work of Choi et al. [16], MonoNHR. It reconstructs a neural radiance field of a clothed human body from a single RGB image without ground truth 3D scans of clothed humans, by conditioning on a human body mesh [34]. While the task and approach are analogous to ours, MonoNHR does not explicitly encode correlation between the object (clothes) and hand (body).

## III. METHOD

The motivation for HandNeRF is to tackle the challenges of 3D scene reconstruction from a 2D RGB image, such as depth uncertainties and partial observation. HandNeRF achieves this by learning hand-object interaction feature that correlates the hand and object geometry, given a 3D hand shape and 2D object segmentation. The overall pipeline of HandNeRF is depicted in Fig. 2. We first elaborate on the theoretical background of HandNeRF and provide the detailed implementation.

### A. Modeling Hand-object Interaction

Consider a point on a 3D object, $\mathbf{x}_o \in \mathbb{R}^3$ where its occupancy or density is $\sigma \in [0, 1]$, i.e., one if occupied, and zero otherwise. The problem of 3D reconstruction of the object can be cast as learning a function that predicts the
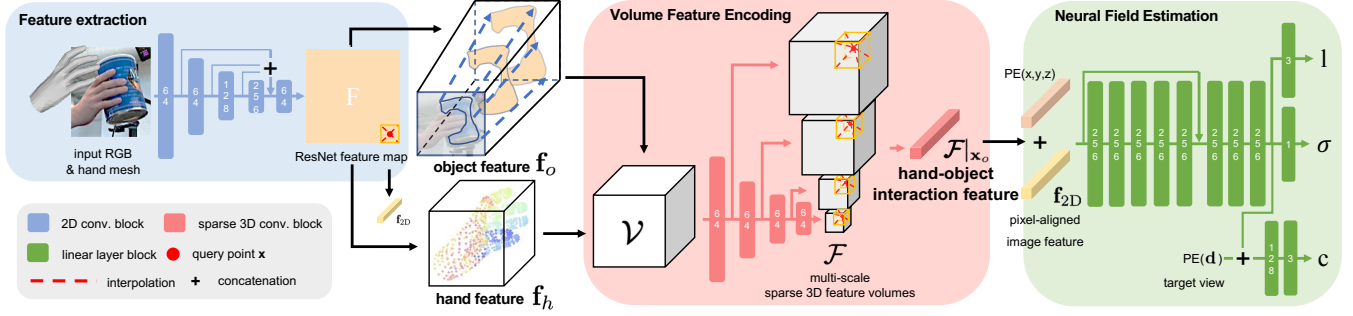
Fig. 2: HandNeRF takes a single RGB image and predicts the volume density, color radiance, and semantic label of each query point in a neural field. Different from comparable works of Ye et al. [15] and Choi et al. [16] that implicitly learns the interaction between hand and object, it explicitly encodes the correlation between hand and object features in 3D space, which provides more accurate 3D reconstruction and novel view synthesis.

density given the location and associated 3D feature $\mathbf{f}_o$:

$$f(\mathbf{x}_o, \mathbf{f}_o) = \sigma, \tag{1}$$

where $f$ is an implicit function of which zero-level set defines the surface of the object. Despite the success of representing objects [14] and humans [35], Equation (1) has a limited capability to express complex scenes such as hand-object interaction scenes.

We extend Equation (1) by incorporating the interactions between the object and hand. Consider a 3D hand mesh $\mathcal{M} = \{\mathbf{m}_i\}$ that is made of a set of faces, where $\mathbf{m}_i$ is the $i^{\text{th}}$ face of the mesh. Each face in the mesh is associated with a 3D feature $\mathbf{f}_h$. We marginalize the density of the object over the density predicted by the hand mesh:

$$f(\mathbf{x}_o, \mathbf{f}_o) = \sum_{\mathbf{x}_h \in \mathcal{M}} f(\mathbf{x}_o, \mathbf{f}_o | \mathbf{x}_h, \mathbf{f}_h) f(\mathbf{x}_h, \mathbf{f}_h), \tag{2}$$

where $\mathbf{x}_h$ is the centroid of the vertices of the face $\mathbf{m}_i$, $f(\mathbf{x}_o, \mathbf{f}_o | \mathbf{x}_h, \mathbf{f}_h)$ is the conditional density given the hand pose and its feature, and $f(\mathbf{x}_h, \mathbf{f}_h) = \{0, 1\}$ is the hand occupancy provided by 3D hand mesh estimation.

Learning $f(\mathbf{x}_o, \mathbf{f}_o | \mathbf{x}_h, \mathbf{f}_h)$ is challenging due to the quadratic complexity of pairwise relationship between all possible pairs of hand and object points $(\mathbf{x}_h, \mathbf{x}_o)$. Instead, we propose to learn an interaction feature $\mathcal{F}$, a correlation between $\mathbf{f}_o$ and $\mathbf{f}_h$ through a series of 3D convolutions:

$$\mathcal{F} = \phi_n \circ \cdots \circ \phi_1 \circ \mathcal{V}, \tag{3}$$

where $\mathcal{F} \in \mathbb{R}^{w \times h \times d \times c}$ is the volume of the interaction features with $w$ width, $h$ height, $d$ depth, $c$ feature dimension, and $\phi_1, \cdots, \phi_n$ are the 3D convolutional filters. The interaction feature $\mathcal{F}|_{\mathbf{x}_o}$ evaluated at an object point $\mathbf{x}_o$ is expected to encode how hand surface points contribute to predicting the occupancy of the point $\mathbf{x}_o$ of the object. The input to the 3D CNN is $\mathcal{V} \in \mathbb{R}^{w \times h \times d \times c'}$, which is the feature volume with the $c'$ feature dimension that includes both hand and object features:

$$\mathcal{V}_{\mathbf{x}} = \begin{cases} \mathbf{f}_h & \text{if} \quad \mathbf{x} \in \{\overline{\mathbf{m}}_i\} \\ \mathbf{f}_o & \text{else if} \quad \Pi\mathbf{x} \in \mathcal{O} \\ \mathbf{0} & \text{otherwise} \end{cases}, \tag{4}$$
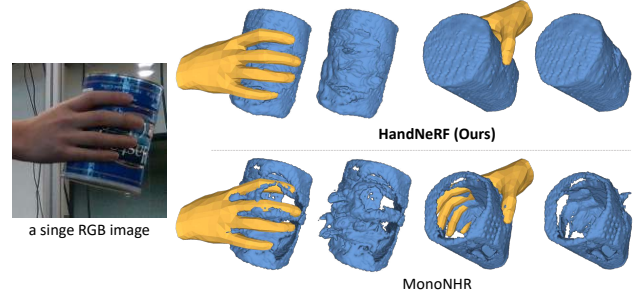


Fig. 3: We visualize object reconstruction with the hand estimation from HandOccNet [5]. Using explicit hand-object interaction features, HandNeRF generates more accurate reconstruction.

where $\mathcal{V}_{\mathbf{x}}$ is the feature at $\mathbf{x}$, $\{\overline{\mathbf{m}}_i\}$ is a set of the centroids of the hand mesh faces, $\Pi\mathbf{x}$ is the camera projection of $\mathbf{x}$ to the input image, and $\mathcal{O}$ is the 2D input object mask.

With the interaction feature $\mathcal{F}$, we extend Equation (1) to include the color, $\mathbf{c} \in \mathbb{R}^3$, and semantic label $\mathbf{l} \in [0, 1]^L$ where $L = 3$ is the number of semantic classes (i.e., hand, object, and background):

$$f(\mathbf{x}_o, \mathbf{d}, \mathbf{f}_{2\mathrm{D}}, \mathcal{F}|_{\mathbf{x}_o}) = (\sigma, \mathbf{c}, \mathbf{l}), \tag{5}$$

where $\mathbf{d}$ is the rendering viewing direction, and $\mathbf{f}_{2\mathrm{D}}$ is the pixel-aligned image feature of $\mathbf{x}_o$. With the prediction of the volume density, color radiance, and semantic label, we render each pixel with its label by integrating the density field. Please refer to Semantic-NeRF [36] for more technical detail of semantic neural rendering.

Fig. 3 shows the impact of the interaction feature $\mathcal{F}$. Our method and MonoNHR [16] both use 3D CNNs to encode features volumetrically based on estimated 3D hand mesh. Unlike MonoNHR, ours explicitly learns hand-object interactions as elaborated, enabling robust object geometry reconstruction even for unobserved and occluded surfaces.

## B. Implementation of HandNeRF

We learn the representation of the hand-object interaction by minimizing the following loss:

$$\mathcal{L} = \sum_{\mathbf{p} \in \mathcal{R}} \left( \left\| \hat{C}(\mathbf{p}) - C(\mathbf{p}) \right\|_2^2 - \sum_{i=1}^{L} L_i(\mathbf{p}) \log(\hat{L}_i(\mathbf{p})) \right)$$

where $\mathcal{R}$ is a set of pixels in multiview images, $\hat{C}(\mathbf{p})$ and $C(\mathbf{p})$ are the predicted and ground truth color of pixel $\mathbf{p}$, respectively, and $\hat{L}_i(\mathbf{p})$ and $L_i(\mathbf{p})$ are the predicted and ground truth semantic label at pixel $\mathbf{p}$.

We design a novel network called *HandNeRF* that predicts a semantic neural radiance field from a single RGB image, as shown in Fig. 2. It is composed of ResNet-18 [37] for feature extraction, sparse 3D convolution layers [38] for volume feature encoding, and linear layers for the neural field estimation. During training, the estimated semantic neural radiance field is validated by projecting on sparse views.

**Input Features:** We deproject a 2D image feature extracted from an image to the points in the 3D volume $\mathcal{V}$ to compose the 3D hand feature $\mathbf{f}_h$ and the 2D object feature $\mathbf{f}_o$ in Equation (4). The 3D hand feature is made of three components using a MANO [21] hand mesh: $\mathbf{f}_h = \begin{bmatrix} \mathbf{h}^\mathsf{T} & \phi(\overline{\mathbf{m}}_i)^\mathsf{T} & \psi(i)^\mathsf{T} \end{bmatrix}^\mathsf{T}$, where $\mathbf{h}$ encodes the visual context of hand-object interaction that is obtained from the 2D image feature at the projection of $\overline{\mathbf{m}}_i$. $\phi(\overline{\mathbf{m}}_i)$ is a positional encoding of the centroid's coordinate, and $\psi(i)$ is the positional encoding of the face index. $\psi(i)$ semantically differentiates a 3D hand point $\mathbf{x}_h$ from different points that are possibly empty, belong to different hand faces, or an object. The 2D object feature is designed in a similar fashion: $\mathbf{f}_o = \begin{bmatrix} \mathbf{o}^\mathsf{T} & \phi(\mathbf{x}_o)^\mathsf{T} & \mathbf{e}^\mathsf{T} \end{bmatrix}^\mathsf{T}$, where $\mathbf{o}$ is the 2D image feature at the projection of $\mathbf{x}_o$, and $\mathbf{e}$ is a constant value for all $\mathbf{x}_o$, where $\mathbf{x}_o \in \{\mathbf{x} \mid \Pi\mathbf{x} \in \mathcal{O} \text{ and } \mathbf{x} \notin \{\overline{\mathbf{m}}_i\}\}$. In practice, $\psi(i)$ and $\mathbf{e}$ are randomly sampled from a Gaussian distribution and fixed during training and testing.

**3D CNN Design:** We correlate the 3D hand feature $\mathbf{f}_h$ and the 2D object feature $\mathbf{f}_o$ with a sparse 3D CNN [38] that takes the feature volume $\mathcal{V}$ as input, to learn the interaction feature. $\mathcal{V}$ rasterizes 3D point coordinates in the neural radiance field with a voxel size of 5mm×5mm×5mm. Before the rasterization, 3D coordinates of object points are perturbated by a random Gaussian noise during training, for augmentation. The sparse 3D CNN produces multi-scale feature volumes, which conceptually add up to the interaction feature volume $\mathcal{F}$ by concatenation along the feature channel dimension. In practice, we keep the feature volumes separated and extract the interaction feature $\mathcal{F}|_{\mathbf{x}}$ of a query point $\mathbf{x}$ per volume with tri-linear interpolation.

## IV. EXPERIMENTS

In this section, we first validate our design of HandNeRF method by conducting detailed ablation studies. Then, we evaluate against the state-of-the-art baselines [14]–[16] that are adapted to be trained on sparse view images without an object template and to be tested on a single image. We adapted IHOI [15] to training with sparse view images,

instead of template-based object annotation as in the original paper, and named it as IHOINeRF. For IHOINeRF, training with semantic labels fails to converge where reconstruction accuracy for hand and object cannot be measured.

**Metrics:** To assess the rendering quality, we use four metrics by comparing with the ground truth images: peak signal-to-noise ratio (PSNR), semantic segmentation intersection over union (IoU), structural similarity index (SSIM), and LPIPS [39]. For the 3D reconstruction accuracy, we compare 3D distance with the ground truth by converting the reconstructed neural radiance field to a 3D mesh using Marching cubes algorithm [40]. F-scores at 5mm and 10mm thresholds, and Chamfer distance (CD) in millimeters are used. We evaluate hand and object separately using 3D segmentation from the predicted semantics.

**Datasets:** We use DexYCB [9] and HO-3D v3 [13] datasets for comparison. In DexYCB, a hand performing object handover is captured from 8 views, and 5 sequences per object are recorded where each sequence shows a distinct grasp pattern. Per object, we keep 4 sequences for training and 1 sequence for testing to validate generalization to novel hand grasps. In Ho-3D v3, an object grasped in a hand is captured from 5 views and 1 sequence per object is recorded where a grasping hand pose changes over time during the sequence. For every object, we split the data to training and testing sets such that the testing set has significantly different grasping hand poses that those in the training set.

## A. Ablation Study

In Table I, we summarize our ablation study to measure the impact of our method design choices. To focus on evaluation of generalization to novel grasp configurations, we train each method per object and average the metrics in the table. We use 4 objects from DexYCB dataset with distinct shapes: 'Cracker Box', 'Banana', 'Power Drill', and 'Coffee Can'. In inference time, all methods use the same 3D hand mesh and 2D object segmentation provided in DexYCB.

**Effect of explicit interaction encoding:** Our main hypothesis is that learning the correlation between hand and object geometry explicitly can regulate the 3D reconstruction of the grasped object given the 3D hand shape. We compare two methods to validate this hypothesis: (M2: $\mathbf{f}_{2D}$) Pixel-NeRF [14] adapted to a single input image that uses the 2D image feature without the 3D hand feature and (M3: $\mathbf{f}_h, \mathbf{f}_o$) a method that uses the 3D hand and 2D object feature without the 2D image feature. As shown in Table I, by exploiting the 3D hand feature, M3 successfully imposes constraints on the relative geometries of hand and object, and provides better generalization to novel hand-object interactions. The results support that our learned interaction feature $\mathcal{F}$ that explicitly encodes hand-object correlations is effective to infer a 3D object geometry without requiring its 3D ground truth during training. The performance of the 3D feature is more pronounced for our method (M5: $\mathbf{f}_h, \mathbf{f}_o + \mathbf{f}_{2D}$) that leverages both the 2D image feature and the 3D feature.

**Significance of 2D object feature:** HandNeRF differs from the existing approaches [15], [16] by using explicit repre-

TABLE I: Ablation study. Our model M5 provides the highest rendering quality, F-scores, and lowest Chamfer Distances (CD) for novel hand-object interaction scenes' object geometry. Subscripts $_w$, $_o$, and $_h$ indicate whole, object, and hand evaluation respectively.

| Method: features | Architect. | PSNR↑ | IoU↑ | SSIM↑ | LPIPS↓ | $F_w$-5 ↑ | $F_w$-10 ↑ | $CD_w$ ↓ | $F_o$-5 ↑ | $F_o$-10 ↑ | $CD_o$ ↓ | $F_h$-5 ↑ | $F_h$-10 ↑ | $CD_h$ ↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| M1: $\mathbf{f}_h$, $\mathbf{f}_o$+$\mathbf{f}_{2D}$ | Transf. | 19.40 | 0.61 | 0.63 | 0.30 | 0.36 | 0.64 | 0.85 | 0.32 | 0.56 | 1.57 | 0.27 | 0.55 | 1.42 |
| M2: $\mathbf{f}_{2D}$ | | 19.09 | 0.61 | 0.60 | 0.31 | 0.30 | 0.56 | 1.17 | 0.28 | 0.47 | 2.90 | 0.19 | 0.49 | 1.56 |
| M3: $\mathbf{f}_h$, $\mathbf{f}_o$ | 3D CNN | 20.11 | 0.77 | 0.65 | 0.27 | 0.47 | 0.78 | 0.30 | 0.41 | 0.68 | 0.62 | **0.54** | **0.92** | 0.12 |
| M4: $\mathbf{f}_h$+$\mathbf{f}_{2D}$ | | 20.31 | 0.72 | 0.68 | 0.26 | 0.40 | 0.68 | 0.79 | 0.30 | 0.53 | 1.73 | **0.54** | **0.92** | **0.09** |
| M5 (ours): $\mathbf{f}_h$, $\mathbf{f}_o$+$\mathbf{f}_{2D}$ | | **21.66** | **0.79** | **0.70** | **0.24** | **0.47** | **0.79** | **0.27** | **0.43** | **0.70** | **0.56** | 0.53 | 0.91 | 0.10 |

TABLE II: Comparison with state-of-the-art baselines on DexYCB and HO-3D v3. Subscripts $_w$, $_o$, and $_h$ indicate whole, object, and hand evaluation, respectively. † indicates use of ground truth 3D hand meshes for inputs, otherwise HandOccNet [5]'s estimation is used. MPJPEs (mean per joint position error) of the estimation are 12mm and 34mm in DexYCB and HO3D v3, respectively.

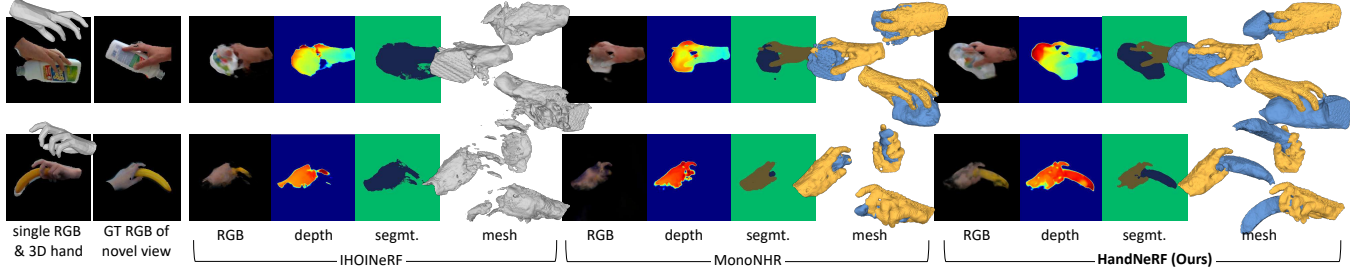| Method | Dataset | PSNR↑ | IoU↑ | SSIM↑ | LPIPS↓ | $F_w$-5 ↑ | $F_w$-10 ↑ | $CD_w$ ↓ | $F_o$-5 ↑ | $F_o$-10 ↑ | $CD_o$ ↓ | $F_h$-5 ↑ | $F_h$-10 ↑ | $CD_h$ ↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PixelNeRF | | 19.09 | 0.61 | 0.60 | 0.31 | 0.30 | 0.56 | 1.17 | 0.28 | 0.47 | 2.90 | 0.19 | 0.49 | 1.56 |
| IHOINeRF | | 18.49 | - | 0.60 | 0.31 | 0.31 | 0.60 | 1.15 | − | − | − | − | − | − |
| IHOINeRF† | | 19.82 | - | 0.64 | 0.27 | 0.38 | 0.69 | 0.54 | − | − | − | − | − | − |
| MonoNHR | DexYCB | 19.36 | 0.63 | 0.64 | 0.30 | 0.37 | 0.62 | 3.19 | 0.25 | 0.43 | 8.59 | 0.49 | 0.89 | 0.11 |
| MonoNHR† | | 19.66 | 0.68 | 0.66 | 0.29 | 0.40 | 0.64 | 3.05 | 0.26 | 0.45 | 8.58 | **0.55** | **0.92** | **0.08** |
| HandNeRF | | 21.19 | 0.75 | 0.68 | 0.25 | 0.44 | 0.77 | 0.31 | 0.42 | 0.68 | 0.59 | 0.46 | 0.88 | 0.12 |
| HandNeRF† | | **21.66** | **0.79** | **0.70** | **0.24** | **0.47** | **0.79** | **0.27** | **0.43** | **0.70** | **0.56** | 0.53 | 0.91 | 0.10 |
| PixelNeRF | | 18.82 | 0.69 | 0.65 | 0.23 | 0.38 | 0.69 | 0.75 | 0.36 | 0.58 | 1.14 | 0.29 | 0.68 | 0.41 |
| IHOINeRF | | 18.55 | - | 0.65 | 0.23 | 0.28 | 0.56 | 1.15 | − | − | − | − | − | − |
| IHOINeRF† | | 19.40 | - | 0.68 | 0.21 | 0.41 | 0.73 | 0.65 | − | − | − | − | − | − |
| MonoNHR | HO3D v3 | 16.98 | 0.66 | 0.60 | 0.21 | 0.28 | 0.53 | 2.09 | 0.19 | 0.37 | 3.49 | 0.33 | 0.65 | 0.64 |
| MonoNHR† | | 19.34 | 0.75 | 0.70 | 0.22 | 0.45 | 0.74 | 0.97 | 0.36 | 0.59 | 1.50 | 0.52 | 0.92 | **0.08** |
| HandNeRF | | 18.04 | 0.68 | 0.63 | 0.23 | 0.38 | 0.70 | 0.35 | 0.40 | 0.66 | 0.41 | 0.45 | 0.51 | 0.78 |
| HandNeRF† | | **20.54** | **0.82** | **0.74** | **0.18** | **0.51** | **0.83** | **0.19** | **0.47** | **0.74** | **0.31** | **0.54** | **0.94** | **0.08** |



Fig. 4: Qualitative results of novel view synthesis (image, depth, and semantic segmentation) and 3D mesh on DexYCB and HO3D v3. Ground truth hand meshes are used as input.

TABLE III: Comparison with state-of-the-art baselines on unseen objects of DexYCB [9]. Subscripts $_w$, $_o$, and $_h$ indicate whole, object, and hand evaluation, respectively. The mesh estimate inputs are from HandOccNet [5]. † indicates using ground truth 3D hand meshes.

| Method | PSNR↑ | IoU↑ | SSIM↑ | LPIPS↓ | $F_w$-5 ↑ | $F_w$-10 ↑ | $CD_w$ ↓ | $F_o$-5 ↑ | $F_o$-10 ↑ | $CD_o$ ↓ | $F_h$-5 ↑ | $F_h$-10 ↑ | $CD_h$ ↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PixelNeRF | 18.85 | 0.58 | 0.58 | 0.34 | 0.21 | 0.47 | 1.09 | 0.22 | 0.44 | 1.31 | 0.10 | 0.31 | 2.09 |
| IHOINeRF | 17.89 | - | 0.58 | 0.34 | 0.21 | 0.45 | 571.55 | − | − | − | − | − | − |
| IHOINeRF† | 19.94 | - | 0.64 | 0.30 | 0.43 | 0.67 | 1.03 | − | − | − | − | − | − |
| MonoNHR | 17.14 | 0.45 | 0.53 | 0.37 | 0.27 | 0.51 | 1.70 | 0.17 | 0.33 | 54.89 | 0.40 | 0.73 | 0.73 |
| MonoNHR† | 19.27 | 0.67 | 0.63 | 0.31 | 0.47 | 0.71 | 1.22 | 0.41 | 0.62 | 1.77 | **0.54** | **0.82** | **0.50** |
| HandNeRF | 18.85 | 0.56 | 0.55 | 0.33 | 0.30 | 0.61 | 0.62 | 0.25 | 0.49 | 0.85 | 0.36 | 0.69 | 0.88 |
| HandNeRF† | **20.83** | **0.72** | **0.66** | **0.27** | **0.51** | **0.75** | **0.72** | **0.46** | **0.68** | **1.11** | 0.51 | 0.80 | 0.52 |

sentation of an object with respect to a 3D hand. To verify the effectiveness of the 2D object feature, we compare two methods: (M4: $\mathbf{f}_h + \mathbf{f}_{2D}$) a method that implicitly learns the hand-object interactions similar to MonoNHR [16] and (M5: $\mathbf{f}_h, \mathbf{f}_o + \mathbf{f}_{2D}$) our method that explicitly models the interactions through the 2D object feature. As shown in Table I, M4 tends to overfit to hand reconstruction and produces poor results for object reconstruction. This implies that without the explicitly defined 2D object feature and its correlation with the hand pose, a strong prior coming from the given hand pose information dominates the prediction while ignoring object information from an input image.
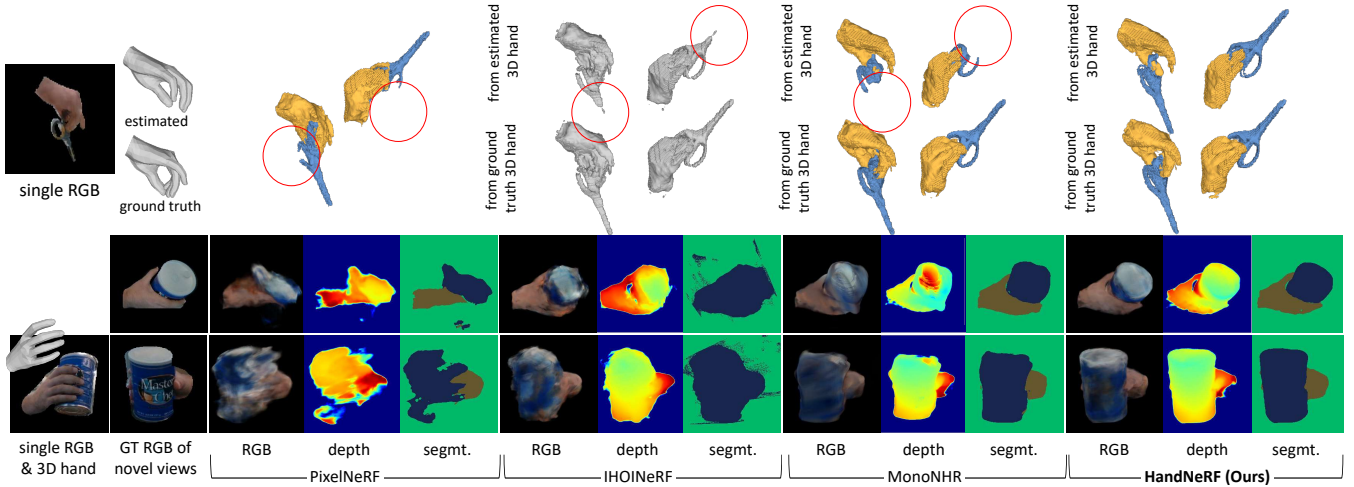
Fig. 5: Qualitative results of novel view synthesis (image, depth, and semantic segmentation) and 3D mesh on DexYCB [9] and HO3D v3 [41], given hand mesh estimation of HandOccNet [5]. The bottom results for scissors are using ground truth hand mesh for reference.
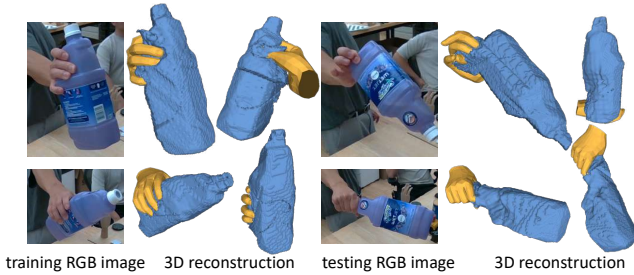


Fig. 6: HandNeRF generalizes to significantly different grasp configurations at the inference time on in-house data. The reconstructed object meshes are visualized with the input hand mesh.
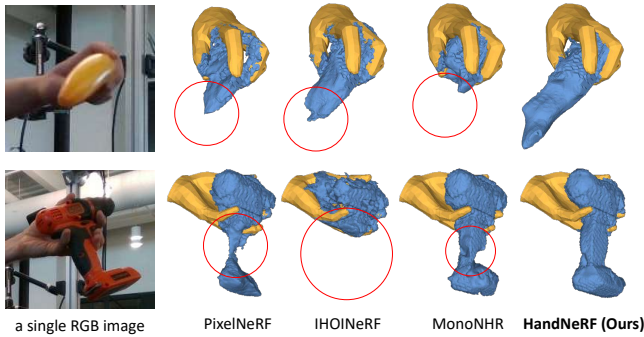


Fig. 7: Qualitative results of generalization to novel object shapes. The reconstructed 3D object mesh is visualized along with the ground truth MANO hand mesh, which is used as input.

**Effect of 3D CNN:** Learning hand-object interaction from Equation (2) is challenging due to the complex quadratic pairwise relationship, requiring a large number of data. Instead, we approximate the quadratic relationship in 3D space using a 3D CNN in Equation (3). We compare against a method (M1) that directly learns the all pairwise relationship between hand and object points using Transformer [43]. As shown in Table I, using the 3D CNN (M5) outperforms M1 in

TABLE IV: Downstream grasping: We compare Contact-GraspNet [42] grasp quality on HandNeRF versus baseline reconstructions for DexYCB handover scenes. HandOccNet hand estimation is used for the methods.

| Input | Reconstruction | Grasp proposal success ratio↑ |
|---|---|---|
| RGB | PixelNeRF | 0.46 |
| | IHOINeRF | 0.42 |
| | MonoNHR | 0.36 |
| | HandNeRF | **0.63** |
| RGBD | - | 0.26 |
| GT mesh | - | 0.77 |

all metrics. Considering the higher gap between training and testing PSNR of M1 (e.g., M1: 4.55 vs. M5: 3.6), the result indicates that our method based on a 3D CNN is resilient to overfitting. Moreover, the model complexity of M1 is over 10 times than ours, comparing the number of model parameters (M1: 27.1K vs. M5: 2.1K).

*B. Comparison with State-of-the-art Methods*

We first evaluate the generalization to novel grasps on the objects seen during training. We assess the generalization to novel object shape by training on 15 DexYCB objects and testing on 4 unseen ones. Finally, we demonstrate the use of reconstruction for grasp planning for robotic handover.

**Generalization to novel grasps:** Table II and Fig. 4 present quantitative and qualitative evaluation on DexYCB and HO3D v3. Given the ground truth grasping hand shape, HandNeRF shows the highest rendering quality scores and the highest F-scores, and the lowest CD (mm) for the whole and object geometry of the novel hand-object interaction scenes. For example, HandNeRF achieves approximately 1.5 times higher F-scores and significantly lower CD for object reconstruction than those of PixelNeRF [14] and MonoNHR [16]. This demonstrates HandNeRF's effective hand-object interaction priors compared to the baselines.

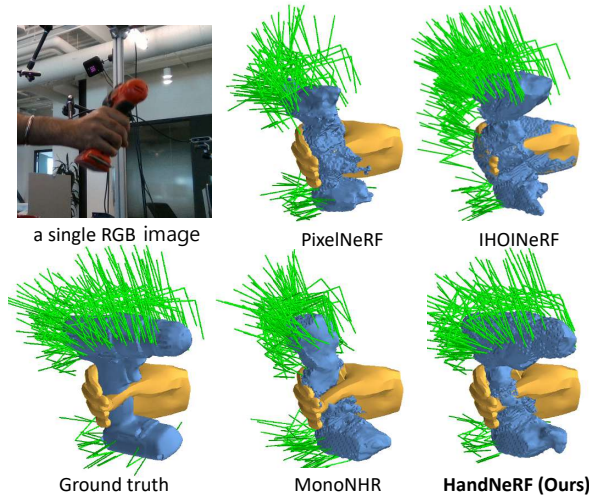We also demonstrate HandNeRF's robustness to erroneous

Fig. 8: Qualitative results of Contact-GraspNet [42]'s grasp proposals on reconstructed meshes of HandNeRF and the baselines. The ground truth MANO hand mesh, which is used as input and grasp collision filtering, is also visualized.

input hand meshes in Fig. 5, which is quantitatively verified in Table II. When using the hand pose estimation from HandOccNet [5], instead of the ground truth, small pose errors in the input hand mesh significantly impact IHOINeRF and MonoNHR outputs. These methods fail to recover half the scissors, implying limitations of implicit interaction encoding. In contrast, given inaccurate 3D hand input, HandNeRF retains reasonable reconstruction quality and renders more accurate novel views far from the input.

We further qualitatively demonstrate HandNeRF's generalization capability on in-house data in Fig. 6. Only 7 RGB cameras are used for data collection and annotation of 3D hand mesh and 2D object segmentation is fully automated. Without ground truth of 3D object geometry during training, HandNeRF exhibits good generalization to significantly different grasping poses and reasonably reconstructs the grasped object from a single RGB image, leveraging the learned hand-object interaction prior.

**Generalization to novel object shape:** As shown in Table III and Fig. 7, HandNeRF consistently outperforms baselines on most metrics, especially reconstructing robust object meshes despite substantial shape dissimilarity between training and test sets. Due to depth ambiguity in a single 2D RGB image, baselines fail to recover overall object shape. For instance, the 'Banana' is only partially reconstructed as its length is unclear from the input. These results demonstrate HandNeRF's superior generalization, likely due to the explicit hand and object geometry encoding effectively regularizing plausible novel object geometry.

**Application to grasp planning for handover:** We evaluate grasp proposals from Contact-GraspNet [42] on reconstructed meshes of HandNeRF and the baselines [14]–[16], RGBD pointcloud of the input image, and ground truth meshes. Grasps colliding with the hand mesh are filtered out before evaluation. We measure the ratio of successful

grasp proposals, where a grasp is counted as successful if it envelops a part of the ground truth object mesh without colliding with it. Unseen scenes of two DexYCB objects ('Banana', 'Power Drill') are used, as Contact-GraspNet performed reliably on their ground truth meshes.

HandNeRF's object reconstruction enables a 1.5 times higher grasp proposal success ratio compared to baselines, as shown in Table IV. Without depth information, HandNeRF achieves a 63% grasp success ratio, approaching the 77% achieved by Contact-GraspNet using ground truth meshes and far exceeding the 26% from the input image pointcloud. Fig. 8 visually demonstrates how HandNeRF's more accurate reconstruction increases successful grasp proposals. Although the surface is locally coarse, HandNeRF's reconstructed global geometry, including the unobserved regions, enables more accurate grasp planning.

## V. LIMITATION AND FUTURE WORK

The limitation of our method in practice is that it strongly depends on hand mesh estimation of off-the-shelf methods. Despite the advance of the recent methods [5], [6], [23], when the hand is severely occluded by the object, the estimated mesh is not accurate enough for inferring further correlation between the hand and object geometry. In such cases, the wrongly estimated hand information can rather heart the object reconstruction. In the future, we will explore to integrate the hand mesh estimation into our system along with the uncertainty modeling to adjust the hand mesh's impact to the final output.

Despite outperformance, our synthesized RGB images are still blurry, when rendered from significantly different view from an input view. Inspired by recent progress on 3D scene generation via language grounding [44], another avenue for future research will be to leverage self-supervised perceptual supervision, such as CLIP [45] feature consistency and object coherency.

## VI. CONCLUSION

This work investigates representation learning for hand-object interactions from a single RGB image. We propose HandNeRF, a method that predicts the semantic neural radiance field of the interaction scenes. The key novelty is the utilization of hand shape to constrain the relative 3D configuration of hands and objects, encoding their correlation explicitly. Unlike existing works, HandNeRF does not require object templates for training and testing, avoiding expensive 3D labeling. Instead, it is supervised with sparse view RGB images, where conventional multi-view reconstruction methods, such as SfM (Structure from Motion), do not apply. HandNeRF outperforms state-of-the-art baselines in rendering and reconstruction on real-world data. Further, we demonstrate improved performance on downstream tasks resulting from HandNeRF's more accurate object meshes, both quantitatively and qualitatively.

## REFERENCES

[1] Bugra Tekin, Federica Bogo, and Marc Pollefeys. H+ o: Unified egocentric recognition of 3d hand-object poses and interactions. In *CVPR*, 2019. 1, 2

[2] Yana Hasson, Bugra Tekin, Federica Bogo, Ivan Laptev, Marc Pollefeys, and Cordelia Schmid. Leveraging photometric consistency over time for sparsely supervised hand-object reconstruction. In *CVPR*, 2020. 1, 2

[3] Shaowei Liu, Hanwen Jiang, Jiarui Xu, Sifei Liu, and Xiaolong Wang. Semi-supervised 3d hand-object poses estimation with interactions in time. In *CVPR*, 2021. 1, 2

[4] Shreyas Hampali, Sayan Deb Sarkar, Mahdi Rad, and Vincent Lepetit. Keypoint transformer: Solving joint identification in challenging hands and object interactions for accurate 3d pose estimation. In *CVPR*, 2022. 1

[5] JoonKyu Park, Yeonguk Oh, Gyeongsik Moon, Hongsuk Choi, and Kyoung Mu Lee. Handoccnet: Occlusion-robust 3d hand mesh estimation network. In *CVPR*, 2022. 1, 3, 5, 6, 7, 11, 12

[6] Yu Rong, Takaaki Shiratori, and Hanbyul Joo. Frankmocap: A monocular 3d whole-body pose estimation system via regression and integration. In *ICCV Workshops*, 2021. 1, 2, 7

[7] Christian Zimmermann, Duygu Ceylan, Jimei Yang, Bryan Russell, Max Argus, and Thomas Brox. Freihand: A dataset for markerless capture of hand pose and shape from single rgb images. In *CVPR*, 2019. 1, 2

[8] Gyeongsik Moon, Shoou-I Yu, He Wen, Takaaki Shiratori, and Kyoung Mu Lee. Interhand2.6m: A dataset and baseline for 3d interacting hand pose estimation from a single rgb image. In *ECCV*, 2020. 1

[9] Yu-Wei Chao, Wei Yang, Yu Xiang, Pavlo Molchanov, Ankur Handa, Jonathan Tremblay, Yashraj S Narang, Karl Van Wyk, Umar Iqbal, Stan Birchfield, et al. Dexycb: A benchmark for capturing hand grasping of objects. In *CVPR*, 2021. 1, 2, 4, 5, 6

[10] Christian Zimmermann, Max Argus, and Thomas Brox. Contrastive representation learning for hand shape estimation. In *GCPR*, 2022. 1, 2

[11] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *CVPR*, 2019. 1, 2, 10

[12] Srinath Sridhar, Franziska Mueller, Michael Zollhöfer, Dan Casas, Antti Oulasvirta, and Christian Theobalt. Real-time joint tracking of a hand manipulating an object from rgb-d input. In *ECCV*, 2016. 1, 2

[13] Shreyas Hampali, Sayan Deb Sarkar, and Vincent Lepetit. Ho-3d_v3: Improving the accuracy of hand-object annotations of the ho-3d dataset. *arXiv*, 2021. 2, 4

[14] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *CVPR*, 2021. 2, 3, 4, 6, 7, 10, 13

[15] Yufei Ye, Abhinav Gupta, and Shubham Tulsiani. What's in your hands? 3d reconstruction of generic objects in hands. In *CVPR*, 2022. 2, 3, 4, 7, 10, 11

[16] Hongsuk Choi, Gyeongsik Moon, Matthieu Armando, Vincent Leroy, Kyoung Mu Lee, and Gregory Rogez. Mononhr: Monocular neural human renderer. In *3DV*, 2022. 2, 3, 4, 5, 6, 7, 13

[17] Bardia Doosti, Shujon Naha, Majid Mirbagheri, and David J Crandall. Hope-net: A graph-based model for hand-object pose estimation. In *CVPR*, 2020. 2

[18] Grégory Rogez, James S Supancic, and Deva Ramanan. Understanding everyday hands in action from rgb-d images. In *ICCV*, 2015. 2

[19] Yana Hasson, Gul Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *CVPR*, 2019. 2, 11

[20] Zhe Cao, Ilija Radosavovic, Angjoo Kanazawa, and Jitendra Malik. Reconstructing hand-object interactions in the wild. In *ICCV*, 2021. 2, 11

[21] Javier Romero, Dimitrios Tzionas, and Michael J Black. Embodied hands: Modeling and capturing hands and bodies together. In *SIGGRAPH Asia*, 2017. 2, 4

[22] Liuhao Ge, Zhou Ren, Yuncheng Li, Zehao Xue, Yingying Wang, Jianfei Cai, and Junsong Yuan. 3d hand shape and pose estimation from a single rgb image. In *CVPR*, 2019. 2

[23] Gyeongsik Moon, Hongsuk Choi, and Kyoung Mu Lee. Accurate 3d hand pose estimation for whole-body 3d human mesh estimation. In *CVPR workshop*, 2022. 2, 7

[24] Korrawe Karunratanakul, Jinlong Yang, Yan Zhang, Michael J Black, Krikamol Muandet, and Siyu Tang. Grasping field: Learning implicit representations for human grasps. In *3DV*, 2020. 2, 11

[25] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan Russell, and Mathieu Aubry. AtlasNet: A Papier-Mâché Approach to Learning 3D Surface Generation. In *CVPR*, 2018. 2

[26] 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 *CVPR*, 2021. 2

[27] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *ICCV*, 2021. 2

[28] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *CVPR*, 2022. 2

[29] Hongyi Xu, Thiemo Alldieck, and Cristian Sminchisescu. H-nerf: Neural radiance fields for rendering and temporal reconstruction of humans in motion. In *NeurIPS*, 2021. 2

[30] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *CVPR*, 2021. 2, 13

[31] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural Human Performer: Learning generalizable radiance fields for human performance rendering. In *NeurIPS*, 2021. 2

[32] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2, 13

[33] Vincent Leroy, Jean-Sébastien Franco, and Edmond Boyer. Volume sweeping: Learning photoconsistency for multi-view shape reconstruction. *IJCV*, 2021. 2

[34] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM TOG*, 2015. 2

[35] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *ICCV*, 2019. 3

[36] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J Davison. In-place scene labelling and understanding with implicit scene representation. In *ICCV*, 2021. 3

[37] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4

[38] Spconv Contributors. Spconv: Spatially sparse convolution library. https://github.com/traveller59/spconv, 2022. 4

[39] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 4

[40] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *SIGGRAPH*, 1987. 4, 13

[41] Shreyas Hampali, Mahdi Rad, Markus Oberweger, and Vincent Lepetit. Honnotate: A method for 3d annotation of hand and object poses. In *CVPR*, 2020. 6

[42] Martin Sundermeyer, Arsalan Mousavian, Rudolph Triebel, and Dieter Fox. Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes. In *ICRA*, 2021. 6, 7

[43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 6

[44] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv*, 2022. 7

[45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 7

[46] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv*, 2021. 10

[47] Zhenggang Tang, Balakumar Sundaralingam, Jonathan Tremblay, Bowen Wen, Ye Yuan, Stephen Tyree, Charles Loop, Alexander Schwing, and Stan Birchfield. Rgb-only reconstruction of tabletop scenes for collision-free manipulator control. In *IEEE IROS*, 2023. 10

[48] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 10

[49] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv*, 2023. 10

[50] Dandan Shan, Jiaqi Geng, Michelle Shu, and David Fouhey. Understanding human hands in contact at internet scale. In *CVPR*, 2020. 10

[51] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 13

[52] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 13

# Supplementary Material of
# HandNeRF: Learning to Reconstruct Hand-Object Interaction Scene from a Single RGB Image

In this supplementary material, we present more experimental results that could not be included in the main manuscript due to the lack of space.

## VII. ADDITIONAL QUALITATIVE RESULTS

We provide more qualitative results and comparison with baselines in Fig. 12 as well as in the online video [1], starting from 01:05 timestamp. In the video, we first show more results of our HandNeRF method for novel view synthesis of RGB, depth, and semantic segmentation, and 3D reconstruction of object and hand meshes. Then, we compare HandNeRF with PixelNeRF [14] and IHOINeRF [15] by rendering RGB and depth images from 360 rotating views, which are significantly different from the input view. Finally, we also demonstrate the effectiveness of HandNeRF's more accurate object reconstruction in downstream tasks such as grasp planning, collision-free motion planning, and object handover. We provide more implementation details on these tasks below.

**Collision-free Motion planning:** In this experiment we demonstrate that accurate object reconstruction is critical to ensure the collision-free motion planning after the object handover to the robot. Specifically, in Isaac Gym simulation [46] environment, we compare the feasibility of the motion plans computed using object reconstructions from PixelNeRF and HandNeRF object. We attach the reconstructed object meshes to the Panda gripper in manually selected grasp configuration. We then compute and save collision-free motions of the robot arm with the attached object reconstruction in a cluttered environment. To evaluate the feasibility of the computed arm motion, we execute the saved robot motion, but with the ground truth object mesh attached to the gripper instead of the reconstruction, similar to [47]. We monitor the motion execution for a potential collision with the environment.

As shown in Fig. 9, we observe that the ground truth object mesh of the power drill collides with the obstacle (a computer display) in the environment when the motion plan using PixelNeRF's object reconstruction is executed. Whereas, the motion plan computed using HandNeRF does not collide with the environment. This motion feasibility verification with the ground truth object mesh demonstrates that the accurate object reconstruction from HandNeRF can enable collision-free motion planning in real world.

**Realworld object handover:** We demonstrate the real world handover results with 'Swiffer Wetjet Refillable Bottle' object which has non-trivial shape. Fig. 10 shows the data collection setup. We use 7 RGB cameras for data collection; 6 RealSense D435 cameras mounted in the scene and 1 RealSense D405 camera attached to the Panda Arm, which
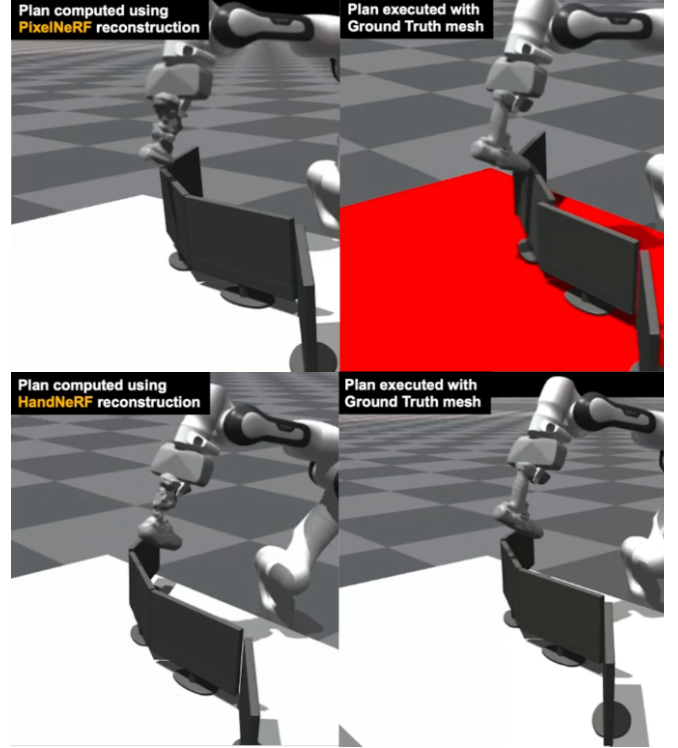


Fig. 9: While HandNeRF's object reconstruction enables successful collision-free motion planning of the Panda Manipulation Arm, inaccurate object reconstruction of PixelNeRF leads to the execution failure. Please refer to the video.

is fixed during the capture. A single computer is used for synchronized capture. Note that our data collection setting is a much more casual setting than conventional multi-view studios that require dozens of synchronized RGBD cameras that entail bandwidth issue for streaming, or involve 3D CAD models of the capturing objects.

Annotation of hand meshes and hand and object segmentation are fully automated. To obtain the 3D hand meshes, we adapted SMPLify-X [11] to MANO hand topology and integrated multi-view 2D projection loss using OpenPose [48] 2D pose detection, which is implemented here [2]. To obtain the segmentation, we used the publically released code of Segment-Anything [49], based on the bounding box detection of Hand Object Detector [50].

The last part of the video shows that 3D reconstruction of HandNeRF enables hand collision-free handover from a single RGB image of novel grasps and object poses in the real world, as shown in Fig. 11. We extended HandNeRF to predict object segmentation for object feature encoding, and the whole process of HandNeRF took approximately 500ms on average. The hand mesh input for HandNeRF is estimated

Fig. 10: Our relatively simple scene capture environment for data collection of hand-object interaction scenes.
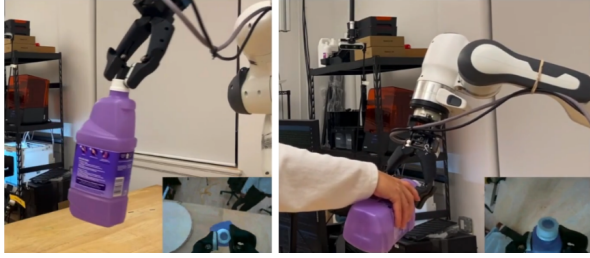


Fig. 11: We demonstrate that the reconstruction of HandNeRF from a single RGB image can used in realworld handover. Please refer to the video.

by HandOccnet [5]. The grasp proposals are from Contact-GraspNet, and we used RRT (Rapidly-exploring random trees) motion planning.

## VIII. GENERALIZATION RESULTS PER OBJECT

We provide quantitative results per object regarding generalization to novel objects, breaking down Table III to Table VII. The split of 15 training objects and 4 testing objects are presented in Table V. We regard the generalization difficulty as 'Power Drill' = 'Banana' > 'Mustard Bottle' > 'Sugar box', considering the difference of shape between training and testing objects.

Table VII shows that HandNeRF generalizes better to more challenging object shapes, that are highly different from training object shapes, in terms of object F-scores and Chamfer distance. For example, MonoNHR and Hand-NeRF are comparable in object reconstruction of 'Sugar Box'. However, HandNeRF outperforms MonoNHR at 62.5 percentage points and 58.1 percentage points for F-score at 10mm threshold, in object reconstruction of 'Banana' and 'Power Drill' respectively. The tendency is similarly observed between PixelNeRF and HandNeRF. The results validate that HandNeRF effectively correlates the possible object geometry to the given hand shape and regularizes it, compared with the baselines.

TABLE V: The training and testing splits of DexYCB, regarding experiments of generalization to novel object in Table III and Table VII.

| Split | Object names |
|---|---|
| training set | 'Coffee Can', 'Cracker Box', 'Tomato Soup Can', 'Tuna Can', 'Pudding Box', 'Gelatin Box', 'Potted Meat Can', 'Pitcher Base', 'Bleach Cleanser', 'Bowl', 'Mug', 'Wood Block', 'Scissors', 'Extra Large Clamp', 'Foam Brick' |
| testing set | 'Sugar Box', 'Mustard Bottle', 'Banana', 'Power Drill' |

TABLE VI: Comparison with the original IHOI [15].

| Method | Hand mesh estimation | | | GT hand mesh | | |
|---|---|---|---|---|---|---|
| | $F_o$-5 ↑ | $F_o$-10 ↑ | $CD_o$ ↓ | $F_o$-5 ↑ | $F_o$-10 ↑ | $CD_o$ ↓ |
| IHOI | 0.40 | 0.67 | 0.62 | 0.65 | 0.84 | 0.25 |
| Ours | 0.40 | 0.66 | 0.41 | 0.54 | 0.74 | 0.31 |

## IX. COMPARISON WITH IHOI

We compare our HandNeRF with the original IHOI [15] that is supervised with 3D object ground truth in Table VI. It outperformed the previous methods [19], [24] that also utilized 3D object ground truth. We test the released model of IHOI on HO3D objects of Table II. The input hand mesh is from HandOccNet [5]'s estimation and object segmentation is from the HO3D v3 annotation.

The model is expected to form an upper performance bound of our HandNeRF as it 1) exploited 3D object ground truth, 2) saw the test scenes' hand grasps and object poses during training from different views, and 3) was pre-trained on MOW dataset [20], while our model did not. Nonetheless, due to our explicit modeling of hand-object interactions, our method produces comparable or better accuracy with the estimated hand mesh from HandOccNet as shown in Table VI.

## X. SENSITIVITY STUDY TO EXTERNAL INPUT

In Table VIII, we conduct additional experiments to evaluate our method's robustness to estimation errors of external inputs, following the work of Ye et al. [15]. For the hand mesh input, we represent the hand pose error of HandOccNet [5] as an estimate deviated from the ground truth annotation. The 100% error source is HandOccNet's predicted hand meshes, the 50% error source is the interpolated hand meshes of the predicted and ground truth hand meshes, and the 200% error source is the extrapolated hand meshes of the predicted and ground truth hand meshes. For the object segmentation input, we evaluate our method with the MaskRCNN segmentation (object mAP: 0.873). Without using any annotation for input, HandNeRF at the bottom row still outperforms the baselines in Table III on DexYCB dataset.

## XI. IMPLEMENTATION DETAIL

All methods in our experimental section follow the same training and testing protocols described below. We will release the full training and testing code along with the data.

TABLE VII: Comparison with the state-of-the-art methods for generalization to novel objects. Per-object evaluation results are presented. The outperformance of HandNeRF becomes more prominent in more challenging objects, i.e. 'Banana' and 'Power Drill', whose shapes are significantly different from those of training objects.

| Method | Object | PSNR↑ | IoU↑ | SSIM↑ | LPIPS↓ | $F_w$-5↑ | $F_w$-10↑ | $CD_w$↓ | $F_o$-5↑ | $F_o$-10↑ | $CD_o$↓ | $F_h$-5↑ | $F_h$-10↑ | $CD_h$↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PixelNeRF | | 16.83 | 0.60 | 0.55 | 0.32 | 0.22 | 0.49 | 0.83 | 0.23 | 0.48 | 0.93 | 0.12 | 0.34 | 1.86 |
| IHOINeRF | | 15.54 | - | 0.52 | 0.36 | 0.26 | 0.52 | 785.03 | — | — | — | — | — | — |
| IHOINeRF† | Sugar | 18.52 | - | 0.63 | 0.27 | 0.42 | 0.67 | 0.55 | — | — | — | — | — | — |
| MonoNHR | Box | 15.58 | 0.47 | 0.48 | 0.39 | 0.29 | 0.56 | 1.17 | 0.20 | 0.40 | 2.23 | 0.44 | **0.77** | 0.76 |
| MonoNHR† | | 17.95 | 0.70 | 0.58 | 0.31 | 0.47 | **0.73** | 0.51 | **0.44** | **0.67** | 0.79 | **0.50** | 0.76 | **0.69** |
| HandNeRF | | 17.29 | 0.57 | 0.51 | 0.33 | 0.32 | 0.62 | 0.53 | 0.25 | 0.50 | **0.75** | 0.41 | 0.74 | 0.85 |
| HandNeRF† | | **19.39** | **0.75** | **0.64** | **0.26** | **0.48** | 0.73 | **0.48** | 0.44 | 0.66 | 0.96 | 0.48 | 0.74 | **0.69** |
| PixelNeRF | | 21.15 | 0.63 | 0.56 | 0.34 | 0.22 | 0.47 | 1.13 | 0.27 | 0.52 | 0.88 | 0.07 | 0.23 | 3.10 |
| IHOINeRF | | 19.39 | - | 0.51 | 0.38 | 0.16 | 0.33 | 1430.62 | — | — | — | — | — | — |
| IHOINeRF† | Mustard | 22.00 | - | 0.60 | 0.31 | 0.37 | 0.62 | 1.49 | — | — | — | — | — | — |
| MonoNHR | Bottle | 17.03 | 0.43 | 0.48 | 0.39 | 0.23 | 0.45 | 2.12 | 0.15 | 0.31 | 3.72 | 0.35 | 0.65 | 1.23 |
| MonoNHR† | | 20.19 | 0.74 | 0.60 | 0.29 | 0.48 | 0.73 | 1.24 | 0.45 | 0.69 | 1.76 | **0.50** | **0.75** | **1.05** |
| HandNeRF | | 20.27 | 0.57 | 0.51 | 0.36 | 0.27 | 0.56 | **0.94** | 0.24 | 0.49 | **1.04** | 0.27 | 0.57 | 1.57 |
| HandNeRF† | | **22.93** | **0.74** | **0.63** | **0.28** | **0.49** | **0.76** | 1.12 | 0.46 | **0.71** | 1.57 | **0.46** | 0.73 | 1.09 |
| PixelNeRF | | 18.88 | 0.54 | 0.66 | 0.32 | 0.20 | 0.44 | 1.44 | 0.17 | 0.32 | 2.23 | 0.12 | 0.36 | 1.54 |
| IHOINeRF | | 19.39 | - | 0.50 | 0.38 | 0.21 | 0.45 | 1.82 | — | — | — | — | — | — |
| IHOINeRF† | | 19.67 | - | 0.70 | 0.28 | 0.48 | 0.71 | 0.77 | — | — | — | — | — | — |
| MonoNHR | Banana | 17.75 | 0.47 | 0.60 | 0.35 | 0.26 | 0.51 | 1.81 | 0.17 | 0.32 | 259.45 | 0.34 | 0.67 | 0.60 |
| MonoNHR† | | 19.45 | 0.61 | 0.72 | 0.28 | 0.47 | 0.66 | 1.77 | 0.36 | 0.51 | 2.45 | **0.58** | **0.88** | **0.13** |
| HandNeRF | | 18.83 | 0.55 | 0.60 | 0.31 | 0.33 | 0.65 | 0.48 | 0.29 | 0.52 | 0.76 | 0.32 | 0.65 | 0.70 |
| HandNeRF† | | **20.59** | **0.70** | **0.75** | **0.25** | **0.54** | **0.76** | **0.38** | **0.49** | **0.67** | **0.52** | 0.56 | 0.87 | 0.14 |
| PixelNeRF | | 19.16 | 0.55 | 0.58 | 0.34 | 0.20 | 0.46 | 1.07 | 0.18 | 0.40 | 1.42 | 0.10 | 0.31 | 1.77 |
| IHOINeRF | | 18.51 | - | 0.57 | 0.35 | 0.22 | 0.47 | 1.90 | — | — | — | — | — | — |
| IHOINeRF† | Power | 19.58 | - | 0.63 | 0.32 | 0.45 | 0.69 | 1.24 | — | — | — | — | — | — |
| MonoNHR | Drill | 18.26 | 0.45 | 0.36 | 0.56 | 0.28 | 0.51 | 1.74 | 0.15 | 0.31 | 3.34 | 0.44 | 0.81 | 0.35 |
| MonoNHR† | | 19.55 | 0.61 | 0.63 | 0.33 | 0.48 | 0.70 | 1.48 | 0.39 | 0.58 | 2.22 | **0.58** | **0.91** | **0.10** |
| HandNeRF | | 19.07 | 0.57 | 0.58 | 0.31 | 0.31 | 0.61 | **0.52** | 0.24 | 0.49 | **0.75** | 0.43 | 0.79 | 0.38 |
| HandNeRF† | | **20.52** | **0.70** | **0.65** | **0.29** | **0.52** | **0.77** | 0.82 | **0.45** | **0.67** | 1.21 | 0.54 | 0.88 | 0.14 |

TABLE VIII: Ablation results of the sensitivity to the external input of HandNeRF. Without using any annotation for input, HandNeRF at the bottom row still outperforms the baselines in Table III on DexYCB dataset.

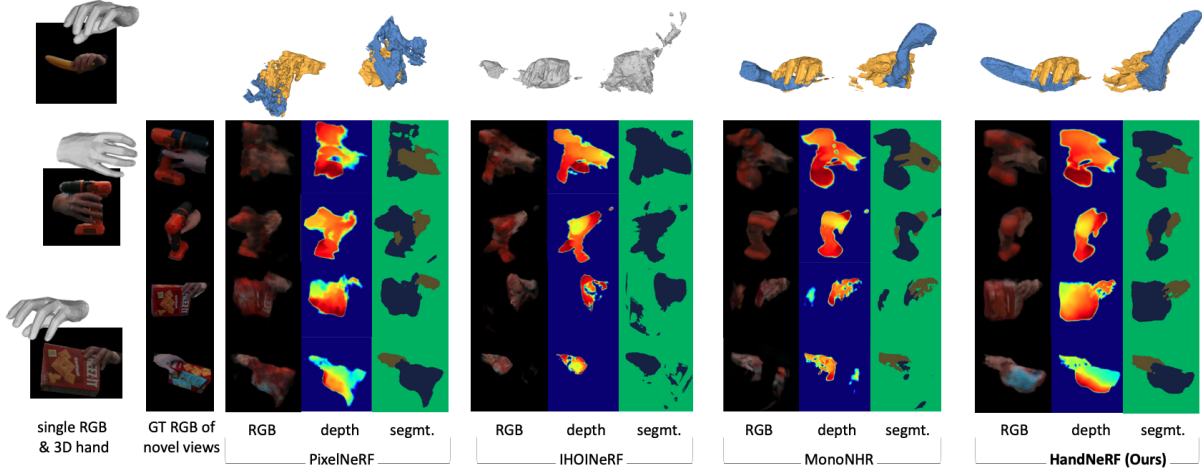| External input type | Source | Dataset | PSNR↑ | IoU↑ | SSIM↑ | LPIPS↓ | $F_w$-5↑ | $F_w$-10↑ | $CD_w$↓ | $F_o$-5↑ | $F_o$-10↑ | $CD_o$↓ | $F_h$-5↑ | $F_h$-10↑ | $CD_h$↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Hand mesh | Annotation | HO3D v3 | 20.54 | 0.82 | 0.74 | 0.18 | 0.51 | 0.83 | 0.19 | 0.47 | 0.74 | 0.31 | 0.54 | 0.94 | 0.08 |
| | 50% error | | 19.08 | 0.74 | 0.67 | 0.20 | 0.44 | 0.79 | 0.23 | 0.43 | 0.70 | 0.36 | 0.44 | 0.86 | 0.15 |
| | 100% error | | 18.04 | 0.68 | 0.63 | 0.23 | 0.38 | 0.70 | 0.35 | 0.40 | 0.66 | 0.41 | 0.45 | 0.51 | 0.78 |
| | 200% error | | 16.81 | 0.59 | 0.58 | 0.28 | 0.30 | 0.59 | 0.95 | 0.36 | 0.61 | 0.50 | 0.17 | 0.38 | 3.22 |
| Segmentation | Annotation | DexYCB | 21.66 | 0.79 | 0.70 | 0.24 | 0.47 | 0.79 | 0.27 | 0.43 | 0.70 | 0.56 | 0.53 | 0.91 | 0.10 |
| | Mask R-CNN | | 20.19 | 0.70 | 0.65 | 0.26 | 0.42 | 0.74 | 0.49 | 0.37 | 0.62 | 0.89 | 0.45 | 0.87 | 0.13 |



Fig. 12: Additional qualitative results. 3D hand mesh estimation of HandOccNet [5] is used as input.

**Training.** We train a method 300 epochs with an initial learning rate of $10^{-3}$. We decay the learning rate by a factor of 10 after 200th epoch. We use the Adam [51] optimizer.

For the volume rendering, we use two separate NeRF modules (i.e., fine and coarse networks) following [14], [32]. We render a one image from a different view per iteration and sample 1024 rays per image. 64 points are sampled along a ray. We gradually increase the ratio of object pixel rays to 0.5 during training, to prevent HandNeRF from being overfitted to the hand rendering. When an object is held by a hand, pixels of hand cover most of the input image area and tend to dominate the volume rendering without this trick. For the feature volume of HandNeRF and MonoNHR [16], we rotate the volume by $\pi/10$ radius around XYZ axes for anti-aliasing of a 3D CNN.

We use Pytorch [52] for code implementation. We use one RTX 3090 GPU during training. It takes approximately 45 hours to fully train HandNeRF, assuming 1000 iterations per epoch.

**Testing.** We perform two steps to evaluate 3D reconstruction. First, we rasterize a hand-object interacting scene with a voxel size of 2mm×2mm×2mm. Then, hand and object meshes are extracted from the estimated volume densities of the neural radiance field using Marching Cubes algorithm [40], following [16], [30]. For networks that estimate a semantic label of a query point $\mathbf{x}$, we separate hand and object voxels before applying the Marching Cubes algorithm to separate hand and object meshes. We remove small meshes with fewer than threshold voxels to reduce the effects of outliers on CD, using connected component graph theory. The threshold is set as a 10% of the number of input voxels.

The reconstruction time of HandNeRF, which includes the process of Marching Cubes algorithm, is 722ms (1016ms), 1015ms (1446ms), 301ms (437ms), 503ms (913ms) per scene for *'002 masterchef can'*, *'003 cracker box'*, *'011 banana'*, *'035 power drill'* of YCB objects respectively. The numbers inside parentheses are reconstruction time with the small mesh sanitization. The rendering time of HandNeRF is 514ms per $3 \times 64 \times 64$ image and 8778ms per $3 \times 256 \times 256$ image.