

ViHOPE: Visuotactile In-Hand Object 6D Pose Estimation with Shape Completion

Hongyu Li^{1,2}, Snehal Dikhale¹, Soshi Iba¹, and Nawid Jamali¹

Abstract—In this letter, we introduce ViHOPE, a novel framework for estimating the 6D pose of an in-hand object using visuotactile perception. Our key insight is that the accuracy of the 6D object pose estimate can be improved by explicitly completing the shape of the object. To this end, we introduce a novel visuotactile shape completion module that uses a conditional Generative Adversarial Network to complete the shape of an in-hand object based on volumetric representation. This approach improves over prior works that directly regress visuotactile observations to a 6D pose. By explicitly completing the shape of the in-hand object and jointly optimizing the shape completion and pose estimation tasks, we improve the accuracy of the 6D object pose estimate. We train and test our model on a synthetic dataset and compare it with the state-of-the-art. In the visuotactile shape completion task, we outperform the state-of-the-art by 265% using the Intersection of Union metric and achieve 88% lower Chamfer Distance. In the visuotactile pose estimation task, we present results that suggest our framework reduces position and angular errors by 35% and 64%, respectively. Furthermore, we ablate our framework to confirm the gain on the 6D object pose estimate from explicitly completing the shape. Ultimately, we show that our framework produces models that are robust to sim-to-real transfer on a real-world robot platform.

Index Terms—Perception for Grasping and Manipulation, Deep Learning for Visual Perception, Force and Tactile Sensing

I. INTRODUCTION

AN accurate 6D pose is a de facto fundamental assumption in numerous applications, such as robotic manipulation [2, 3], autonomous driving [4], and social navigation [5]. The absence of precise knowledge of the pose of the object makes it challenging for an agent to interact with it accurately or avoid it effectively. Recently, deep learning approaches have demonstrated promising results [6–8]. These methods, when combined with iterative refinement [6, 9, 10], leverage the object’s 3D model to obtain a more accurate estimate. However, many methods do not perform well in the presence of intermediate to extreme occlusions, particularly in scenarios involving dexterous manipulation where the object is being held, grasped, or sometimes completely obscured by the robot hand. In such scenarios, finding an accurate pose of the partially observed shape is challenging.

Manuscript received: April 18, 2023; Revised July 19, 2023; Accepted August 25, 2023.

This paper was recommended for publication by Editor Markus Vincze upon evaluation of the Associate Editor and Reviewers’ comments.

¹The authors are with Honda Research Institute USA, Inc. {snehalsubhash_dikhale, siba, njamali}@honda-ri.com

²Hongyu Li is with Brown University hongyu@brown.edu. This work was completed when he was an intern at Honda Research Institute USA, Inc. Digital Object Identifier (DOI): see top of this page.

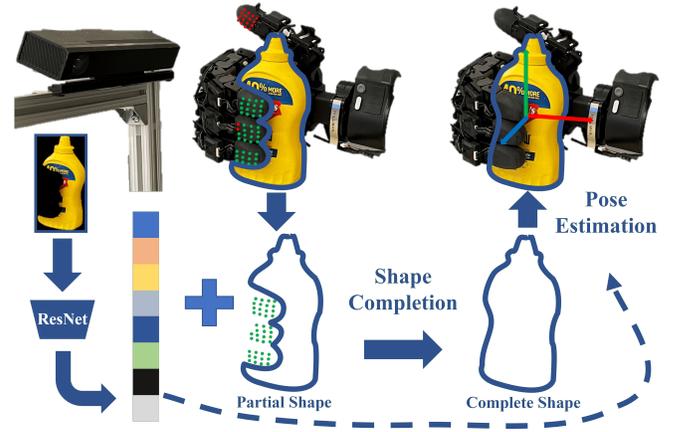


Fig. 1: A high-level overview of the proposed framework. The green dots represent taxels of the tactile sensors that are in contact. RGB image and depth map are retrieved from an RGB-D sensor for object segmentation and visual feature extraction using ResNet [1]. The visual features, tactile data, and partial shape observation (from the vision sensor) are fed into our shape completion module and transformed into a complete shape of the in-hand object. Finally, the completed shape (in latent space) is used to estimate the object’s 6D pose through a joint optimization of both shape and pose.

Consequently, researchers have considered inclusion of tactile sensors into the sensor suite to improve the quality of perception during manipulation [11–13]. Villalonga et al. [13] leverage a template-based approach to match the visuotactile observation with the rendered shapes. Dikhale et al. [12] use visuotactile data and utilize an end-to-end deep neural network and address the 6D pose estimation as a regression problem. However, they do not explicitly leverage the 3D geometry of the object.

To this end, we introduce ViHOPE—Visuotactile In-Hand Object Pose Estimator (Fig. 1). ViHOPE takes visuotactile observation as input and explicitly optimizes the object shape while estimating the pose. We hypothesize that jointly optimizing the shape and pose of the object will provide more accurate estimates of the 6D pose of the object. Additionally, during deployment, by providing an estimate of the complete shape of the in-hand object, it increases explainability and potentially expands the range of applications, such as grasping [14]. Specifically, we first train an autoencoder to capture the geometry prior of the object and encode the object shape into a latent space. We then leverage a GAN to transfer the latent

code from the partial shape space to that of the complete shape. We then use the estimated complete latent code, and the visual feature to estimate the 6D pose.

We conduct experiments on a synthetic dataset by Dikhale et al. [12] and a physical robot platform. We evaluate the performance of ViHOPE on two tasks: shape completion and pose estimation. In the shape completion task, we show improved performance by a large margin compared with the prior work [15], where our model faithfully reconstructs the complete shape even under heavy occlusion. In the pose estimation task, we demonstrate our model outperforms the state-of-the-art visuotactile pose estimator [12]. We also present results of ablation studies, in which we remove the shape completion module to confirm the effectiveness and robustness of our approach, which explicitly optimizes shape.

In summary: 1) We propose a novel visuotactile shape completion module based on the volumetric representation that accurately recovers the complete shape of the object under heavy occlusion. 2) We present a novel framework for visuotactile 6D object pose estimation that optimizes the shape completion and the pose estimation modules jointly, leveraging object geometry to improve the pose estimation.

In this letter, for simplicity, unless otherwise indicated, we use the phrase *pose estimation* for visuotactile instance-level in-hand 6D object pose estimation and the phrase *shape completion* for visuotactile shape completion.

II. RELATED WORKS

In this section, we briefly summarize the related literature from two aspects: shape completion and 6D pose estimation.

A. Shape Completion

The objective of the shape completion task is to estimate the complete shape of an object from a partial observation. Wu et al. [16] approach the shape completion problem using cGAN. Zhang et al. [17] formulate the problem using GAN inversion. However, their method has a time complexity that is 3,500 times greater than direct methods [18], rendering it impractical for real-time robotic applications.

In robotic manipulation research, to counter self-occlusion and occlusions in a cluttered environment, prior works utilize shape completion in a modular [14], and in an end-to-end manner [19]. Wang et al. [20] reconstruct object shape using RGB image followed by shape refinement using tactile data. Watkins-Valls et al. [15] complete the object's shape using visuotactile data using CNNs; however, their approach is not validated on in-hand objects.

B. 6D Pose Estimation

In the field of 6D pose estimation, the problem is approached by combining various modalities, including but not limited to vision, and tactile. In this section, we will present the literature on this subject by categorizing it into two broad categories: non-visuotactile-based and visuotactile-based.

Non-visuotactile-based: The task of 6D object pose estimation has been centered around geometry. One popular trend

is to find the correspondence between the observation and the 3D model of the object, such as 2D-3D correspondence [7, 21–23] using PnP/RANSAC [24] or 3D-3D correspondence [10] using registration methods like ICP [9]. In these works, the exact 3D model of the object is provided during both the training and inference time. Park et al. [21] and Wang et al. [22] detect the 2D pixel-wise normalized coordinate map of the object and match it with the object's 3D model using PnP/RANSAC. While these approaches exhibit impressive performance in estimating pose on novel instances, the non-differentiable correspondence matching prohibits end-to-end training with downstream tasks. To counter this limitation, later works [7, 23] investigate differentiable correspondence matching and learn the pose in an end-to-end manner.

On the other hand, numerous works use the end-to-end method to estimate the 6D pose without explicitly finding correspondence. One of the pioneering works, PoseCNN [25] estimates the 6D pose from monocular RGB input. Gao et al. [26] estimates 6D pose using only point cloud, leveraging PointNet [27] to extract geometry features from point cloud. DenseFusion [6] utilizes both RGB and point cloud data and fuse the extracted features [1, 27] on a pixel-to-pixel basis to estimate the pose. The above-mentioned approaches demonstrate promising improvement on established benchmarks like LINEMOD [28], or YCB-Video [25]. However, their performance has not been evaluated under heavy occlusion typically found in in-hand manipulation.

Visuotactile-based: Several studies utilize both visual and tactile modalities to address the challenge of occlusions in in-hand objects. Villalonga et al. [13] employ a template-based approach to estimate the pose of in-hand objects by rendering a set of shape images from a 3D model in various random poses, then the recorded shape by the tactile sensor is compared with the collection to determine the most probable pose. Dikhale et al. [12] regress the RGB-D and tactile data to a 6D pose in an end-to-end manner and outperform the prior RGB [25] and RGB-D [6] approaches. However, they do not explicitly leverage the geometry of the object.

Unlike prior works, ViHOPE takes advantage of both the end-to-end approach and the object geometry. We focus on jointly optimizing the object shape and the pose estimation to improve the accuracy of the pose estimate.

III. METHODOLOGY

Given an image \mathcal{I} and its respective depth map \mathcal{D} , an object \mathcal{O} , and the tactile feedback \mathcal{T} from the robot hand, our goal is to estimate the 6D pose $[R|t]$ of the object \mathcal{O} with respect to the camera, where R represents the 3D rotation, and t represents the 3D translation. We assume the 3D model of the object \mathcal{O} is available during training.

Unlike vision sensors, tactile sensors are less developed [29]. There is no standard representation format for tactile data. To ensure the robustness of pose estimation algorithms against variations in tactile sensor types, we adopt the approach proposed by Dikhale et al. [12], in which the tactile feedback \mathcal{T} is presented to the algorithm in the form of object-surface point cloud $P^{\mathcal{T}}$. That is, when the tactile

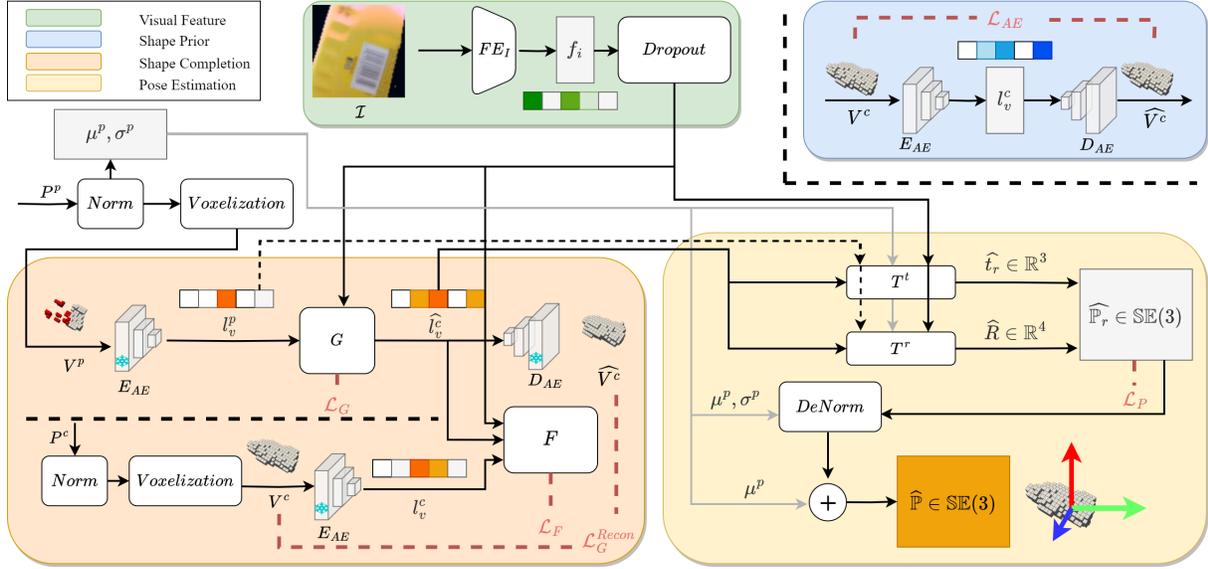


Fig. 2: The proposed framework consists of three phases. Phase one trains the autoencoder (the blue box) in isolation to learn a shape prior. Phase two trains the shape completion module (the orange box) with frozen weights for E_{AE} and D_{AE} from phase one. Phase three uses the completed shape (in latent space) to estimate the object’s pose in an end-to-end manner.

sensors make contact with the object, the robot’s kinematic is used to get the 3D position of each taxel in contact with the object, using a user-defined threshold on the sensor feedback.

The proposed visuotactile pose estimation framework (Fig. 2) consists of: a visuotactile shape completion module and a pose estimation module. The shape completion module explicitly optimizes the object’s shape from a partial observation of the object from the vision and tactile sensors (Section III-A). The pose estimation module uses the output of the shape completion module, in the form of latent code, along with the visual features and point cloud normalization scalars to estimate the 6D pose of the object (Section III-B).

A. Shape Completion

Our framework starts from a volumetric shape completion module, which consists of two steps: 1) learning a shape prior from the object model, i.e., full observation, 2) recovering the complete shape by learning a mapping from partial observation to a complete one, in the latent space.

1) *Learning the shape prior:* We first train an autoencoder (the blue box in Fig. 2) to capture the shape prior of the object. To capture the prior under different orientations, we apply random $SO(3)$ rotations to the object \mathcal{O} , and voxelize it to form an augmented dataset. The autoencoder consists of two components: an encoder E_{AE} and a decoder D_{AE} . The encoder encodes the input voxel occupancy grid $V \in \mathbb{R}^{n_x \times n_y \times n_z}$ into a latent code $l_v \in \mathbb{R}^{n_l}$ using a set of 3D convolutional layers. The decoder recovers the original voxel occupancy grid from the latent code as $\hat{V} \in \mathbb{R}^{n_x \times n_y \times n_z}$ using symmetrical 3D deconvolutional layers. We apply a batch normalization layer after each layer, followed by a ReLU activation function. We set $n_x = n_y = n_z = 32$ and $n_l = 128$,

empirically, and optimize the autoencoder model using the Jaccard index loss

$$\mathcal{L}_{AE} = 1 - \frac{|V^c \cap D_{AE}(E_{AE}(V^c))|}{|V^c \cup D_{AE}(E_{AE}(V^c))|}. \quad (1)$$

Note, higher voxel grid resolution can improve accuracy but comes with increased computation and memory costs. Some works have explored methods to mitigate these costs [30, 31]; however, such optimization techniques are beyond the scope of this work.

2) *Recovering the complete shape:* After training the autoencoder to convergence, we optimize the entire shape completion model (the orange box in Fig. 2). The encoder E_{AE} and the decoder D_{AE} of the shape completion module are initialized with the weights from the previous step and are frozen in this step [16].

The inputs to the shape completion model consist of occluded observational data, that is, a partial visuotactile point cloud $P^p \in \mathbb{R}^{3 \times N}$ and an RGB image $\mathcal{I} \in \mathbb{R}^{H \times W \times 3}$. We first obtain a semantic segmentation mask \mathcal{S} of the object \mathcal{O} using the RGB image \mathcal{I} . We then segment the depth map \mathcal{D} using the mask \mathcal{S} as \mathcal{D}^p and transform \mathcal{D}^p to its respective point cloud form P^D . Combining the point cloud P^D with the tactile point cloud P^T , we obtain the observation $P^p = P^D \cup P^T$ of the object \mathcal{O} . To improve the training efficiency, we first normalize the input partial point cloud P^p using its centroid $\mu^p \in \mathbb{R}^3$ and the farthest distance from the centroid $\sigma^p \in \mathbb{R}$. The normalized point cloud is voxelized as V^p and encoded using the frozen encoder E_{AE} into a partial latent vector l_v^p in the partial latent space \mathcal{M}^p . During training, we use the ground-truth complete voxel grid V^c to calculate the losses.

To recover the complete shape from the partial observation, we seek a mapping for shape latent code from the partial latent space to the complete latent space $\mathcal{M}^p \mapsto \mathcal{M}^c$. Inspired by

Wu et al. [16], we find the mapping using a cGAN [32]. We condition the cGAN on the partial latent vector l_v^p and the visual feature f_i , which is extracted from the input image \mathcal{I} using a pretrained ResNet [1] feature extractor FE_I . FE_I is finetuned during the training period and regularized by a dropout layer. The dropout layer is only activated during training and deactivated during testing. Therefore our generator is trained as $G : (\mathcal{M}^p, p(f_i)) \mapsto \mathcal{M}^c$. We pass the estimated complete latent vector $\hat{l}_v^c \in \mathcal{M}^c$ from the generator to the discriminator F along with the ground-truth complete latent vector l_v^c , obtained by applying the same procedure on the ground-truth complete point cloud P^c . Having the same conditioning as G , the discriminator F is trained as a binary classifier to distinguish the real complete latent vector l_v^c and the fake complete latent vector \hat{l}_v^c . At the end, we feed the estimated complete latent vector \hat{l}_v^c into the frozen decoder D_{AE} to reconstruct the complete shape \widehat{V}^c .

The shape completion model is optimized using three losses: the discriminator loss \mathcal{L}_F , the generator loss \mathcal{L}_G , and the reconstruction loss \mathcal{L}_G^{Recon} . The discriminator loss penalizes the discriminator if it can't distinguish the real and fake latent vectors. On the other side of this min-max game, the generator loss encourages the generator to fool the discriminator. We leverage the loss functions from LSGAN [33] to stabilize the training

$$\mathcal{L}_F = \mathbb{E}_{V^c, \mathcal{I}} [F_{cGAN}(E_{AE}(V^c), FE_I(\mathcal{I})) - 1]^2 + \mathbb{E}_{V^p, \mathcal{I}} [F_{cGAN}(G_{cGAN}(E_{AE}(V^p), FE_I(\mathcal{I}))) - 1]^2 \quad (2)$$

$$\mathcal{L}_G = \mathbb{E}_{V^p, \mathcal{I}} [F_{cGAN}(G_{cGAN}(E_{AE}(V^p), FE_I(\mathcal{I}))) - 1]^2. \quad (3)$$

To further stabilize the GAN training and guide the model, we include the reconstruction loss \mathcal{L}_G^{Recon} that directly measures the differences between the ground-truth complete shape V^c and the estimated complete shape $\widehat{V}^c \triangleq D_{AE}(G_{cGAN}(E_{AE}(V^p), FE_I(\mathcal{I})))$ using Jaccard index loss

$$\mathcal{L}_G^{Recon} = 1 - \frac{|V^c \cap \widehat{V}^c|}{|V^c \cup \widehat{V}^c|}. \quad (4)$$

Therefore, the overall training objective for our shape completion model is

$$\operatorname{argmin}_{G, FE_I} \operatorname{argmax}_F \mathcal{L}_F + \mathcal{L}_G + \alpha \mathcal{L}_G^{Recon}, \quad (5)$$

where α is the weight for the reconstruction loss.

B. Pose Estimator

We use two simple four-layer MLP models to estimate the pose of the object. Our pose estimators T^t and T^r take the estimated complete shape latent code \hat{l}_v^c , visual features f_i , normalization factors μ^p and σ^p , and a skip connection from partial latent l_v^p as input and estimate the 3D translation residual $\hat{t}_r \in \mathbb{R}^3$ and 3D rotation in quaternion form $\widehat{R} \in \mathbb{R}^4$, respectively. Note, similar to the prior works [6, 12], instead of estimating the absolute translation t , we estimate the residual

of the translation $t_r = t - \mu^p$. We use the residual pose $\widehat{\mathbb{P}}_r = [\widehat{R} | \hat{t}_r]$ to calculate the point cloud loss \mathcal{L}_P [6]:

$$\mathcal{L}_P = \frac{1}{k} \sum_{x \in \mathcal{K}} \|(Rx + t_r) - (\widehat{R}x + \hat{t}_r)\|, \quad (6)$$

where \mathcal{K} denotes a set of points randomly sampled from the object's 3D model, and k represents the cardinality $|\mathcal{K}|$. The point cloud loss minimizes the distance between the points on the ground-truth pose and their respective points on the models transformed using the estimated pose. The overall loss function is shown in Equation 7,

$$\operatorname{argmin}_{G, FE_I, T^t, T^r} \operatorname{argmax}_F \mathcal{L}_F + \mathcal{L}_G + \alpha \mathcal{L}_G^{Recon} + \beta \mathcal{L}_P, \quad (7)$$

where β is the weight for the point cloud loss.

To speed up convergence, our method is trained in four steps: i) The autoencoder ($E_{AE} + D_{AE}$) is trained in isolation, and weights are frozen; ii) The shape completion module ($G + F$) is trained; iii) The pose estimator ($T^t + T^r$) is trained while freezing the shape completion module; iv) the shape completion module in unfrozen and trained end-to-end ($G + F + T^t + T^r$), similar to [34, 35].

IV. EXPERIMENTS

Our experiments are designed to assess the efficacy of our proposed framework: 1) in the level of accuracy achieved by our shape completion module in reconstructing object shapes; 2) the impact of explicit shape completion on the quality of 6D pose estimation; 3) the contribution of each component of the framework to pose estimation accuracy; and 4) the sensitivity of the framework to variations in occlusion levels and tactile contact points. This section outlines our experimental setup including model training and performance evaluation. The model was trained and tested on a synthetic dataset, and then transferred to a real physical robot to study the framework's robustness in sim-to-real transfers.

Synthetic Dataset: We use VisuoTactile synthetic dataset from Dikhale et al. [12] to train our framework. In this dataset, a subset of 11 YCB objects [36] are selected based on their graspability. A total number of 20K distinct in-hand poses are simulated per object. In particular, Unreal Engine 4.0 has been used to render photo-realistic observational data of a 6 DoF robot arm with a 4-fingered gripper equipped with 12 32x32 tactile sensors (3 per finger). A main RGB-D camera captures images of the robot holding an object. Each tactile sensor captures object surface contact points in a point cloud format. Each data sample is generated by randomizing the in-hand object pose, the robot fingers configuration, and the robot arm orientation and position. Domain randomization is also applied for the color and pattern of the background and workspace desk.

Real Robot: We test our model on the Allegro Hand (Wonik Robotics) and the Sawyer robot arm (Rethink Robotics). Our Allegro Hand is instrumented with the uSkin 4x4 tactile sensors and uSkin Curved tactile sensors from XELA Robotics. Each finger has three 4x4 tactile sensors (16 taxels) and one Curved tactile sensor (30 taxels). We use a Microsoft Kinect V2 camera to collect RGB-D data.

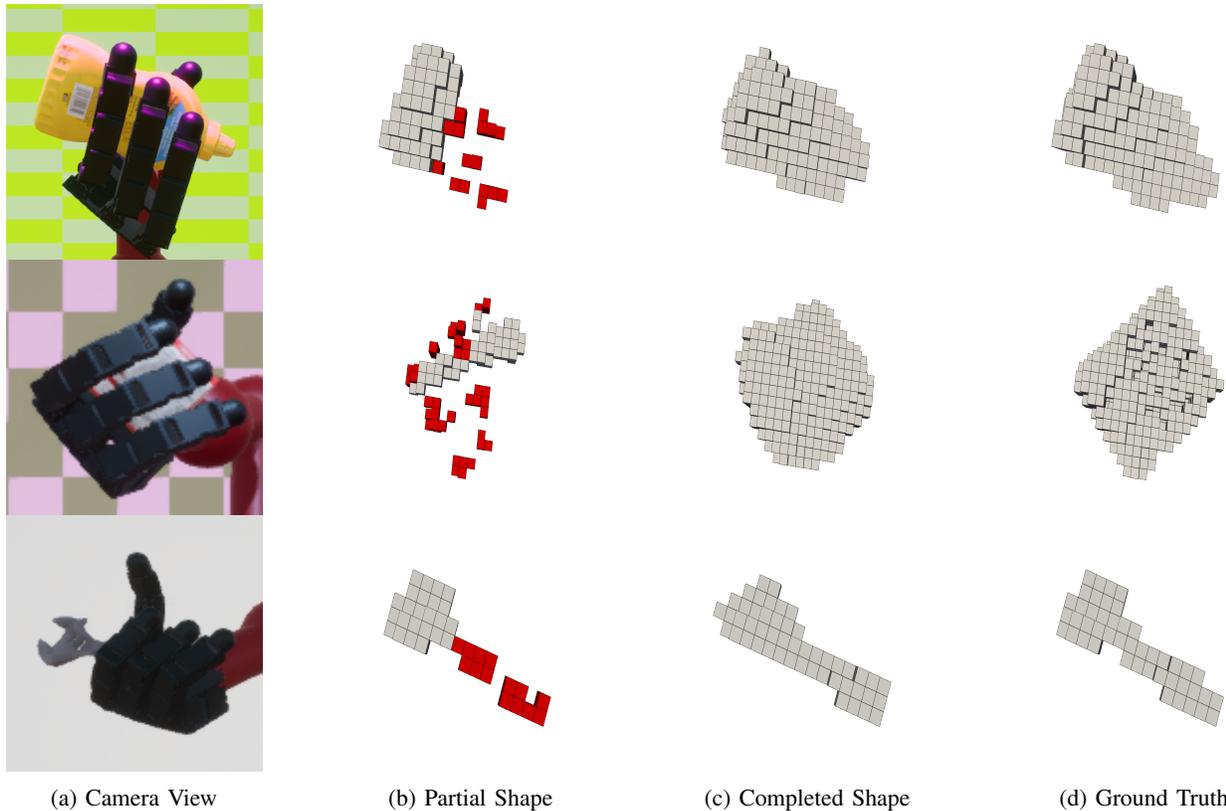


Fig. 3: Visuotactile shape completion results. The gray and red voxel represent RGB-D and tactile observations, P^D and P^T , respectively. From left to right, we show the cropped RGB image from the main camera, the observed partial shape of the in-hand object, the completed shape by our approach, and the ground-truth shape. From top to bottom are the samples drawn from the YCB categories “006_mustard_bottle”, “009_gelatin_box”, and “044_adjustable_wrench”.

Implementation Details: We utilize an ImageNet pre-trained ResNet-34 model as our visual feature extractor FE_I . We apply center crop, Gaussian blur, and color jitter augmentations to the input image. We implement our pose estimators T^r and T^t using four-layer MLPs with layers of [512, 256, 32, 4] and [512, 256, 32, 3], respectively. We trained our autoencoder for 500 epochs, shape completion module for 1000 epochs, pose estimator for 1000 epochs, and the entire end-to-end model for 2000 epochs using the Adam optimizer with a learning rate of 1×10^{-3} , 5×10^{-4} , 1×10^{-3} , and 5×10^{-5} , respectively. We set the reconstruction loss α and the point cloud loss β (Equation 7) to 30 and 5000, respectively.

V. RESULTS

A. Shape Completion

1) *Qualitative Result:* Fig. 3 represents a visualization of the input and output of our shape completion module, which uses visuotactile observations to faithfully complete the shape of the in-hand object. As the 3D reconstruction can have varying levels of occlusion from different viewing angles, we recommend viewing the accompanying video for a clearer and improved visualization.

2) *Quantitative Result:* We compare our shape completion module with the seminal work from Watkins-Valls et al. [15] using two metrics: Intersection over Union (IoU) and Chamfer

TABLE I: Quantitative shape completion result.

Method	IoU \uparrow	CD \downarrow
Watkins-Valls et al. [15]	0.142	0.125
Vision Only	0.341	0.042
ViHOPE (Ours)	0.519	0.015

Distance (CD). We implement their proposed model [14, 15] using PyTorch with a minor modification. To ensure a fair evaluation, a consistent voxel occupancy grid resolution was used. Our implementation utilizes a 32^3 voxel grid for input and output, instead of the 40^3 voxel grid utilized in the original work. This modification allowed seamless integration of the shape completion module into our established pipeline. Results in Table I show a 265.5% increase in IoU and an 88% decrease in CD, demonstrating the robustness of our model under challenging conditions.

We also evaluate the performance of the shape completion model by removing the tactile modality (Vision Only). The result suggests the significance of tactile modality for completing the shape of an in-hand object.

B. 6D Pose Estimation

The Metrics: We evaluate the performance of our pose estimator using two metrics: position error and angular error. The position error is determined as the L2 norm of the

TABLE II: A comparison of our approach with the state-of-the-art is presented in the first half of the table, followed by the results of our ablation studies in the second half. The modalities used by the methods are highlighted in the second column.

Method	Modalities			Position Error (cm) ↓	Angular Error (deg) ↓	ADD (cm) ↓	ADD-S (cm) ↓
	RGB	point cloud	tactile				
PoseCNN [25]	✓			6.146 ± 0.023	10.897 ± 0.082	-	-
DenseFusion [6]	✓	✓		0.640 ± 0.004	9.969 ± 0.117	1.037 ± 0.008	0.571 ± 0.003
ViTa [12]	✓	✓	✓	0.299 ± 0.002	8.074 ± 0.105	0.825 ± 0.007	0.474 ± 0.002
No-Shape-Completion	✓	✓	✓	0.258 ± 0.018	4.104 ± 0.049	0.400 ± 0.019	0.282 ± 0.018
No-Vis-GAN	✓	✓	✓	1.613 ± 0.333	4.132 ± 0.058	1.745 ± 0.333	1.623 ± 0.332
No-Vis-MLP	✓	✓	✓	0.156 ± 0.001	5.677 ± 0.083	0.403 ± 0.004	0.214 ± 0.001
No-Tactile	✓	✓		1.614 ± 0.015	17.228 ± 0.165	2.023 ± 0.017	0.774 ± 0.009
No-Point-Cloud	✓		✓	0.861 ± 0.010	14.245 ± 0.096	1.478 ± 0.011	0.655 ± 0.009
ViHOPE (Ours)	✓	✓	✓	0.194 ± 0.009	2.873 ± 0.036	0.298 ± 0.009	0.214 ± 0.008

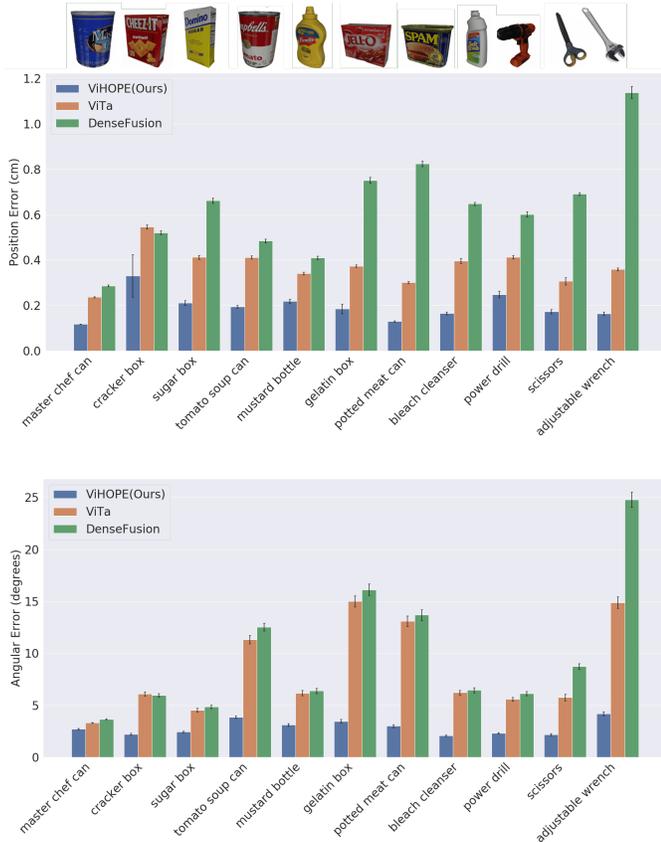


Fig. 4: Performance comparison with the state-of-the-art.

difference between the estimated translation vector and the ground truth translation vector, $\|t - \hat{t}\|_2$. The angular error is computed as the inverse cosine of the inner product of the estimated rotation quaternion and the ground truth quaternion, $\cos^{-1}(2\langle R, \hat{R} \rangle - 1)$.

1) *Comparison with state-of-the-art*: We compare the performance of our pose estimation network with two seminal works: i) the visuotactile-based estimator (ViTa) [12], and ii) the RGB-D-based estimator, DenseFusion [6]. In Fig. 4, we provide a per-instance numerical analysis on 11 YCB objects. Our approach outperforms ViTa and DenseFusion by a large margin on each object with statistical significance, suggesting explicit shape optimization is more effective compared to

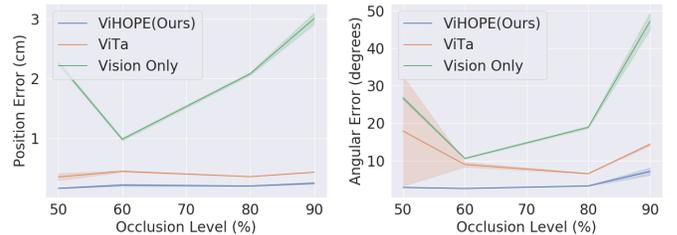


Fig. 5: Performance under different levels of occlusion

implicit methods.

2) *Performance under different occlusion level*: We evaluate the performance of our model under different levels of occlusion. The results of our evaluation are presented in Fig. 5, where it can be observed that the model demonstrates a robust performance in the presence of increasing levels of occlusion. It is worth noting that our method maintains its performance as compared to ViTa, which suggests that our model is able to handle occlusion effectively and still produce competitive performance. We further evaluate the performance by removing the tactile modality (Vision Only). The result confirms the significant contribution of tactile modality under severe occlusions.

3) *Performance under different tactile points*: The spatial resolution of real-world tactile sensors can vary significantly. Therefore, we analyze the performance of our model under different tactile contact points (Fig. 6). It is noteworthy that our model, which was trained with 80 tactile points, demonstrates robust performance when presented with a reduced number of points. As expected, the performance of the model deteriorates as the number of tactile points is reduced and drops significantly when the tactile modality is removed entirely. It is worth noting that compared to ViTa, which requires tactile input, our model can still provide pose estimation even without tactile feedback, although with degraded performance. Upon analyzing the position error, we observe that, up to a reduction of 40 tactile points, our model outperforms ViTa, which uses 1000 points. The angular error results show that our model consistently outperforms ViTa.

4) *Ablation Studies*: We perform ablation studies to examine the effectiveness of our design choices (Table II).

No-Shape-Completion: To evaluate the contribution of explicit shape completion, we remove the shape completion

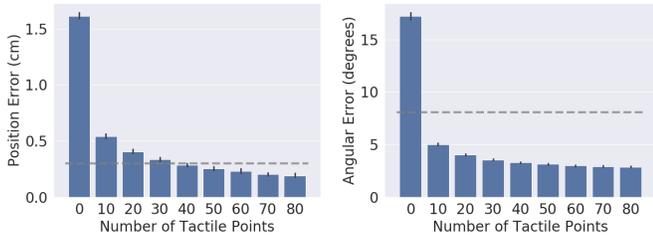


Fig. 6: The performance of our approach under different levels of tactile contact points. The dashed gray line represents the performance of ViTa using 1000 tactile contact points.

module from our proposed framework, which is achieved by feeding the partial latent vector l_v^p to the pose estimator instead of the complete vector \hat{l}_v^c . Our ablation result shows that by removing the shape completion module, the position error drops by 24.8%, and the angular error drops by 30.0%, suggesting, jointly and explicitly optimizing the shape and the pose during visuotactile pose estimation is effective.

No-Vis-Gan: To analyze the gain from visual cues in the shape completion module, we remove the visual feature conditioning f_i . We observe that removing the visual conditioning from the shape completion module resulted in a significant deterioration of performance, highlighting the importance of incorporating visual cues. The study shows that, without visual cues, the partial geometry feature is ambiguous for inferring the complete shape under heavy hand occlusion.

No-Vis-MLP: To examine the contribution of visual features in pose estimators, we remove the visual feature input from the pose estimators. We notice that removing the visual features degrades the angular error performance. This makes sense because our dataset contains symmetrical objects. The object geometry alone is insufficient for accurately determining the pose of symmetrical objects, such as a mustard bottle, which requires the utilization of visual features to distinguish between its front and back.

No-Tactile & No-Point-Cloud: Two separate studies were conducted to evaluate the contribution of the tactile points P^T (No-Tactile) and the point cloud from the vision sensor P^D (No-Point-Cloud). Our results suggest a significant drop in performance when either the tactile points or the point cloud input from the visual sensors is removed, emphasizing the significant contribution of both modalities to pose estimation.

5) *Real-world experiment*: We validate the sim-to-real robustness of our framework using our robot platform, with a subset of YCB objects that could be grasped by the Allegro Hand. The hand moves along a trajectory covering different poses. We apply a novel hand-grasping pose that doesn't exist in our training dataset, which shows our model's ability to generalize. Due to the excessive point cloud noise from the RGB-D sensor in the real world, we apply a point cloud statistical outlier filter as pre-processing using Open3D [37]. We consider 20 neighbors with a standard deviation ratio of 1. Our model efficiently operates at 11.2ms / 89Hz using an NVIDIA RTX 6000, thus capable of real-time deployment.

In Fig. 7, we show three consecutive frames of the Allegro

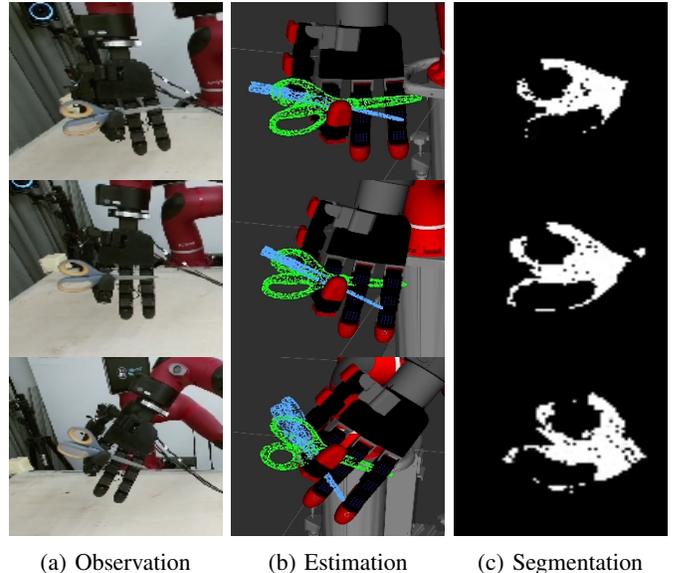


Fig. 7: We compare our method (in green) against ViTa [12] (in blue) in real-world. From left to right are the observation from the RGB-D sensor, the pose estimations, and the noisy partial segmentation task we obtain.

Hand holding a scissor. Our model could accurately estimate the 6D pose of the in-hand object, given a noisy and partial segmentation mask. In the bottom row, we demonstrate a failure case of our method where the estimated pose of the scissors is flipped 180 degrees due to the near-symmetrical geometry and visual feature of the scissors. We refer readers to our supplementary videos for our experiment videos that include more objects and real-time quantitative comparisons.

VI. CONCLUSION

We presented ViHOPE, a novel framework for estimating the 6D pose of an in-hand object using visuotactile perception. We introduced explicit shape completion, which we hypothesize improves the accuracy of the 6D pose estimate by jointly optimizing both the shape and pose of the object. We validated our hypothesis by conducting experiments using a synthetic dataset of 11 YCB objects and compared ViHOPE's performance with state-of-the-art methods. Our results suggest a 35% reduction in position error and a 64% reduction in angular error in the pose estimation task compared to the state-of-the-art. We also demonstrated that ViHOPE outperforms state-of-the-art shape completion approaches by 265% in terms of IoU and 88% lower in CD. We presented the results of ablation studies that confirmed the contribution of explicit shape completion to the accuracy of the 6D pose estimate. To assess the practical viability of our framework in situations where access to high-resolution tactile sensors may be limited, we conducted experiments evaluating its performance under reduced tactile contact points. Our findings indicate that our framework outperforms the current state-of-the-art and can still produce reasonable pose estimates even in the absence of tactile feedback, although with a decreased performance. Finally, we validated the sim-to-real robustness of ViHOPE

in a real-world robot experiment, suggesting its success in estimating the 6D pose of an object in real-world settings.

In our study, we used an empirically determined voxel resolution of 32^3 . Using a coarser grid would compromise performance, while a finer grid would introduce computational overhead and latency. The voxel resolution is crucial in capturing the level of detail in the object's shape. Higher resolutions capture finer geometric features. This aspect becomes particularly important for objects whose accurate pose determination depends on distinguishing geometric features smaller than the selected voxel grid resolution. Moving forward, we plan to explore octree representations [30, 31] that enable higher resolutions while maintaining computational efficiency.

We focus on instance-level pose estimation in this letter. In the future, we are interested in extending our pose estimator to work on more challenging scenarios, for example, lack of annotated data [38], and category-level pose estimation [22]. Another interesting future direction is incorporating other sensory information, such as pressure, into our framework. Another direction that can be pursued is the use of methods such as long short-term memory networks to use temporal coherence to filter out erroneous estimates.

REFERENCES

- [1] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," 2016.
- [2] Y. Qin, H. Su, and X. Wang, "From One Hand to Multiple Hands: Imitation Learning for Dexterous Manipulation from Single-Camera Teleoperation," Apr. 2022, arXiv:2204.12490 [cs].
- [3] T. Chen, J. Xu, and P. Agrawal, "A System for General In-Hand Object Re-Orienting," in *Proceedings of the 5th Conference on Robot Learning*. PMLR, Jan. 2022.
- [4] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? The KITTI vision benchmark suite," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2012.
- [5] Y. F. Chen, M. Everett, M. Liu, and J. P. How, "Socially aware motion planning with deep reinforcement learning," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*.
- [6] C. Wang, D. Xu, Y. Zhu, R. Martin-Martin, C. Lu, L. Fei-Fei, and S. Savarese, "DenseFusion: 6D Object Pose Estimation by Iterative Dense Fusion," 2019.
- [7] H. Chen, P. Wang, F. Wang, W. Tian, L. Xiong, and H. Li, "EPro-PnP: Generalized End-to-End Probabilistic Perspective-N-Points for Monocular Object Pose Estimation," 2022.
- [8] Y. Hu, J. Hugonot, P. Fua, and M. Salzmann, "Segmentation-Driven 6D Object Pose Estimation," 2019.
- [9] P. J. Besl and N. D. McKay, "Method for registration of 3-D shapes," in *Sensor Fusion IV: Control Paradigms and Data Structures*, vol. 1611. SPIE, Apr. 1992.
- [10] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox, "DeepIM: Deep Iterative Matching for 6D Pose Estimation," 2018.
- [11] J. Bimbo, S. Luo, K. Althoefer, and H. Liu, "In-Hand Object Pose Estimation Using Covariance-Based Tactile To Geometry Matching," *IEEE Robotics and Automation Letters*, vol. 1, no. 1, Jan. 2016.
- [12] S. Dikhale, K. Patel, D. Dhingra, I. Naramura, A. Hayashi, S. Iba, and N. Jamali, "VisuoTactile 6D Pose Estimation of an In-Hand Object Using Vision and Tactile Sensor Data," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, Apr. 2022.
- [13] M. B. Villalonga, A. Rodriguez, B. Lim, E. Valls, and T. Sechopoulos, "Tactile Object Pose Estimation from the First Touch with Geometric Contact Rendering," in *Proceedings of the 2020 Conference on Robot Learning*. PMLR, Oct. 2021.
- [14] J. Varley, C. DeChant, A. Richardson, J. Ruales, and P. Allen, "Shape completion enabled robotic grasping," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sep. 2017.
- [15] D. Watkins-Valls, J. Varley, and P. Allen, "Multi-Modal Geometric Learning for Grasping and Manipulation," in *2019 International Conference on Robotics and Automation (ICRA)*, May 2019.
- [16] R. Wu, X. Chen, Y. Zhuang, and B. Chen, "Multimodal Shape Completion via Conditional Generative Adversarial Networks," in *Computer Vision – ECCV 2020*, ser. Lecture Notes in Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Cham: Springer International Publishing, 2020.
- [17] J. Zhang, X. Chen, Z. Cai, L. Pan, H. Zhao, S. Yi, C. K. Yeo, B. Dai, and C. C. Loy, "Unsupervised 3D Shape Completion Through GAN Inversion," 2021.
- [18] Y. Cai, K.-Y. Lin, C. Zhang, Q. Wang, X. Wang, and H. Li, "Learning a Structured Latent Space for Unsupervised Point Cloud Completion," 2022.
- [19] Z. Jiang, Y. Zhu, M. Svetlik, K. Fang, and Y. Zhu, "Synergies Between Affordance and Geometry: 6-DoF Grasp Detection via Implicit Representations," Jul. 2021, arXiv:2104.01542 [cs].
- [20] S. Wang, J. Wu, X. Sun, W. Yuan, W. T. Freeman, J. B. Tenenbaum, and E. H. Adelson, "3D Shape Perception from Monocular Vision, Touch, and Shape Priors," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2018.
- [21] K. Park, T. Patten, and M. Vincze, "Pix2Pose: Pixel-Wise Coordinate Regression of Objects for 6D Pose Estimation," 2019.
- [22] H. Wang, S. Sridhar, J. Huang, J. Valentin, S. Song, and L. J. Guibas, "Normalized Object Coordinate Space for Category-Level 6D Object Pose and Size Estimation," 2019.
- [23] G. Wang, F. Manhardt, F. Tombari, and X. Ji, "GDR-Net: Geometry-Guided Direct Regression Network for Monocular 6D Object Pose Estimation," 2021.
- [24] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, Jun. 1981.
- [25] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "PoseCNN: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes," May 2018, arXiv:1711.00199 [cs].
- [26] G. Gao, M. Lauri, Y. Wang, X. Hu, J. Zhang, and S. Frintrop, "6D Object Pose Regression via Supervised Learning on Point Clouds," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, May 2020.
- [27] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017.
- [28] S. Hinterstoisser, S. Holzer, C. Cagniard, S. Ilic, K. Konolige, N. Navab, and V. Lepetit, "Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes," in *2011 International Conference on Computer Vision*, Nov. 2011.
- [29] A. Yamaguchi and C. G. Atkeson, "Recent progress in tactile sensing and sensors for robotic manipulation: can we turn tactile sensing into vision?" *Advanced Robotics*, vol. 33, no. 14, Jul. 2019.
- [30] M. Tatarchenko, A. Dosovitskiy, and T. Brox, "Octree Generating Networks: Efficient Convolutional Architectures for High-Resolution 3D Outputs," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct. 2017.
- [31] H. Li, Z. Li, N. U. Akmandor, H. Jiang, Y. Wang, and T. Padir, "Stereo-voxelnet: Real-time obstacle detection based on occupancy voxels from a stereo camera using deep neural networks," in *2023 International Conference on Robotics and Automation (ICRA)*.
- [32] M. Mirza and S. Osindero, "Conditional Generative Adversarial Nets," Nov. 2014, arXiv:1411.1784 [cs, stat].
- [33] X. Mao, Q. Li, H. Xie, R. Y. K. Lau, Z. Wang, and S. P. Smolley, "Least Squares Generative Adversarial Networks," Apr. 2017, arXiv:1611.04076 [cs].
- [34] Y. Liang, B. Chen, and S. Song, "SSCNav: Confidence-Aware Semantic Scene Completion for Visual Semantic Navigation," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*.
- [35] T. Gervet, S. Chintala, D. Batra, J. Malik, and D. S. Chaplot, "Navigating to Objects in the Real World," Dec. 2022, arXiv:2212.00922 [cs].
- [36] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar, "The YCB object and Model set: Towards common benchmarks for manipulation research," in *2015 International Conference on Advanced Robotics (ICAR)*, Jul. 2015.
- [37] Q.-Y. Zhou, J. Park, and V. Koltun, "Open3D: A Modern Library for 3D Data Processing," Jan. 2018, arXiv:1801.09847 [cs].
- [38] K. Zhang, Y. Fu, S. Borse, H. Cai, F. Porikli, and X. Wang, "Self-Supervised Geometric Correspondence for Category-Level 6D Object Pose Estimation in the Wild," Oct. 2022, arXiv:2210.07199 [cs].