# DONet: Learning Category-Level 6D Object Pose and Size Estimation from Depth Observation

Haitao Lin    Zichang Liu    Chilam Cheang    Lingwei Zhang    Yanwei Fu    Xiangyang Xue

Fudan University

## Abstract

*We propose a method of Category-level 6D Object Pose and Size Estimation (COPSE) from a single depth image, without external pose-annotated real-world training data. While previous works [43, 39, 4] exploit visual cues in RGB(D) images, our method makes inferences based on the rich geometric information of the object in the depth channel alone. Essentially, our framework explores such geometric information by learning the unified 3D Orientation-Consistent Representations (3D-OCR) module, and further enforced by the property of Geometry-constrained Reflection Symmetry (GeoReS) module. The magnitude information of object size and the center point is finally estimated by Mirror-Paired Dimensional Estimation (MPDE) module. Extensive experiments on the category-level NOCS benchmark demonstrate that our framework competes with state-of-the-art approaches that require labeled real-world images. We also deploy our approach to a physical Baxter robot to perform manipulation tasks on unseen but category-known instances, and the results further validate the efficacy of our proposed model.*

## 1. Introduction

This paper studies the task of the "Category-level" 6D Object Pose and Size Estimation (COPSE), which aims at addressing the limitation of previous "instance-level" object pose estimation [28, 42, 12, 26, 19] that requires exact object 3D models. Typically, estimating accurate 6D object pose plays a pivotal role in the tasks of augmented reality [23], scene understanding [37], and robotic grasping and manipulation [6, 40, 5, 22] depicted in Fig. 1.

However, learning COPSE is fundamentally very challenging, due to the large intra-category variation of object categories, and the huge domain gap between real and synthetic training images. We explain these challenges next.
(1) The *intra-class variation* is caused by color/texture and geometry/shape discrepancies of various objects within a
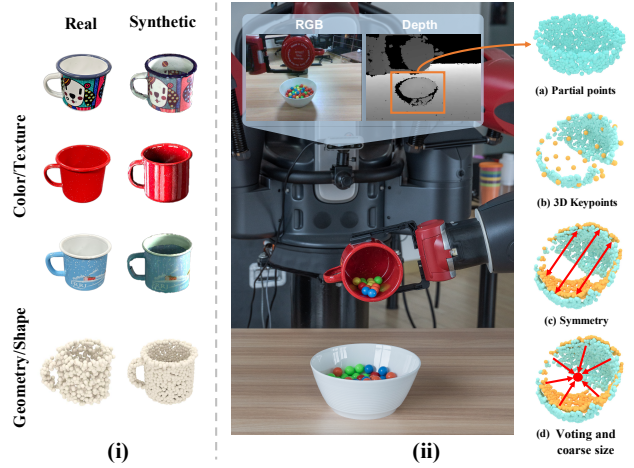


Figure 1: (i) The comparisons between the texture/color and the shape/geometry under the synthetic and real-world image domain. (ii) Deploy our DONet in a physical Baxter robot to execute the pouring task. Given (a) partial points (back-projected from the masked depth map), our point-based network reasons (b) semantically aligned keypoints (yellow points) via the 3D-OCR module, (c) symmetry correspondences via the GeoRes module, and (d) dense voting vectors and rough object size via the MPDE module. Then 6D object pose and size estimation is implemented.

category [32, 33, 34]. To address this challenge, most previous works resort to utilizing the unified space, e.g., Normalized Object Coordinate Space (NOCS) [43] and Canonical Shape Space (CASS) [4]. Essentially, as the empirical evaluation in [39, 43], the performance of previous works solely rely on the color information in RGB channels. Hence, it remains an open question how the geometry information from the depth channel should be best exploited. Take the `mug` in Fig. 1(i) as an example. We show the mug of the same shape, and yet different colors. Various cup colors of RGB information cause visual diversity while all cups geometrically have nearly the same shape in depth channel. This motivates us to better exploit geometry from the depth channel to alleviate the intra-class variation.

(2) Typically, due to limited pose-annotated real-world images, most COPSE models are trained on synthetic data and tested by real images, respectively. This leads to the *synthetic-real image domain gap*. Such a gap reflects the different distribution of synthetic and real data [43]. As shown in Fig. 1(i), such domain gap negatively affects more on RGB channels than the depth channel, due to different lighting conditions and material textures. Furthermore, This also hinders real-world COPSE applications, as plenty of synthesized training examples cannot be effectively leveraged to real testing images.

In general, our recipe is to better exploit depth information for the COPSE model, to improve its generalization ability. Essentially, it is nontrivial to employ depth information to alleviate the intra-class variation and synthetic-real domain gap, in contrast to the previous "densefusion" manner [42] of directly extending "instance-level" RGB-D feature fusion methods to category-level tasks. Critically, the design principle of our model comes from the following two points.

*Shape-correspondence-oriented orientation modeling.* Assume the category of the observed object is known, the predefined template point cloud of the corresponding category is deformed to align with the observed object point cloud from the depth channel. Such deformed template points are dubbed as an implicit representation of object orientation, computed by the 3D-OCR module. The representation in Fig. 1(ii)(b) (yellow points) has both a good template shape and semantically well aligned to the observable points in Fig. 1(ii)(a), thus reducing the intra-category variation. Such representation facilitates the calculation of orientation, according to visually semantic shape consistency of category-specific instances. Theoretically, 3D-OCR learns shape correspondence or variation between template shape and intra-class instances. The nature of 3D-OCR should alleviate intra-class variation, thus allows it more generalized in handling non-linear rotation space.

*Symmetry-oriented modeling.* We reason the potential mirrored points for the observable points. We illustrate this by yellow points in Fig. 1(c). Particularly, we enforce the property of geometry-constrained reflection symmetry by the GeoReS module. The magnitude information of object size and center point is further estimated by Mirror-Paired Dimensional Estimation (MPDE) module, as in Fig. 1(d). Note that, GeoRes utilizes symmetry reasoning to help object pose and size estimation. For asymmetric objects, the symmetry principle should be still usable by GeoRes to make a rough estimation of object pose and size.

Formally, this paper proposes a geometry-based approach for category-level 6D object pose and size recovery, from depth information represented by point clouds. Concretely, our method has three well-designed modules, including: (1) **3D-OCR**, which extracts the point-wise embedding vector features to effectively capture the intra-class shape correspondence information for the semantically shape-guided reconstruction. (2) **GeoReS** estimates the mirrored points of the input observable points, validated to encourage the overall performance. Also, learning GeoReS is a much easier task than direct object completion which normally demands a larger and more complicated network. (3) **MPDE** decouples the dimensional information of object center and size from the roughly complete shape. Extensive experiments conducted on category-level datasets demonstrate that our method outperforms the state-of-the-art by using synthetic data only, with significant performance improvement ($5°2cm$ mAP: 24.1% vs. 19.3%, and $5°5cm$ Accuracy: 49.1% vs. 33.3%).

To sum up, our work makes the following contributions:
1) A novel intermediate 3D-OCR is proposed in category-level object orientation estimation, which implicitly learns semantic shape alignment of the observable point cloud for rotation recovery.
2) We present a novel GeoReS component to leverage the reflection symmetry assumption as an auxiliary enhancement for pose and size estimation.
3) A MPDE component is proposed for accurate object center and size estimation. And jointly training modules aforementioned further improves the performance.
4) We give the implementation of the COPSE model on the real Baxter robot. It shows that our 6D object pose and size estimator achieves nearly real-time performance, running at around 10Hz on a single NVIDIA RTX2070 GPU.

## 2. Related Work

**Instance-Level methods.** Most previous works [21, 11, 18, 2] only estimate instance-level object poses, by directly matching the image features. Unfortunately, these methods are less efficient to infer the pose of texture-less objects. Recent efforts are made on directly regressing 6D object pose from RGB images by CNN-based architectures, such as PoseNet [15] and PoseCNN [46]. DenseFusion [42] introduces "densefusion" manner for better aggregating color and depth information from RGB-D images, which infers more accurate objects pose than RGB-only methods. Such fusion manner is also used in recent category-level tasks [4, 39]. Another line of works [20, 25, 27, 28, 31, 38, 48] firstly regress the object coordinates or keypoints in 2D images, and then recover poses by Perspective-n-Point algorithm [17], e.g., PVNet [28]. Heuristically, PVN3D [10] expends keypoints voting mechanism from 2D to 3D space. Different from PVNet and PVN3D methods, our approach focuses on a more general setting without exact object 3D models in practical applications.

**Category-level methods.** Recent COPSE approaches vitally alleviate the limitation of previous instance-level tasks.
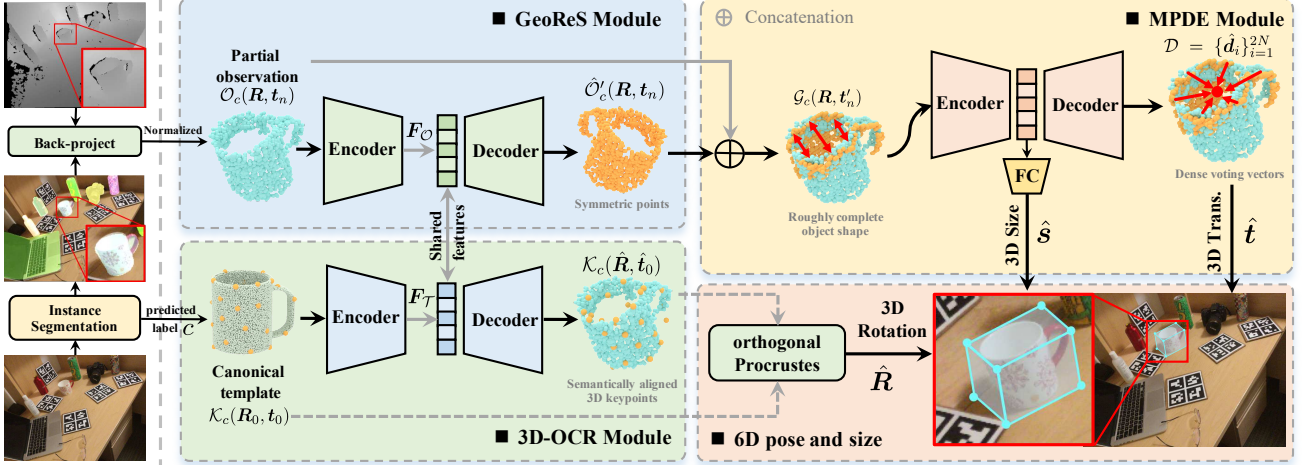
Figure 2: Architecture overview. The pre-processing stage (left) outputs the predicted category labels and potential masks of target instances (e.g., mug). The back-projected points from the masked depth map and canonical category-specific template points are sent to our network (right). The network is composed of four output branches that generate 3D keypoints, symmetric points, dense voting vectors, and normalized scale. Then the final 6D object pose and size are recovered by the post-processing step.

Sahin *et al.* [34] first address COPSE problem by training a part-based random forest with 3D skeleton for pose recovery. Wang *et al.* [43] propose a method to learn pixel-wise correspondence, mapping foreground image pixels to 3D object coordinates in NOCS. Chen *et al.* [4] leverages pose-independent 3D shape representation for complete model reconstruction in CASS. Tian *et al.* [39] model object shape via coarse-to-fine deformation upon the prior mean shape, and learns correspondence for mapping points of the reconstructed model into unified space, i.e., NOCS. 6-PACK [41] is the first category-specific pose tracker to achieve real-time tracking that assumes available initial object pose. On the contrary, our pose estimator capable of *multi-category* estimation has a more lightweight model for efficient pose and size inference than the pose tracking method.

**Symmetry.** Symmetry property has been widely adopted in recent works [7, 24, 44]. The reconstruction of symmetry objects has been investigated in [47, 14]. Wu *et al.* [45] use latent symmetry properties to disentangle components obtained from a single image. In the field of 6D pose estimation, HybridPose [36] is the first work to take the dense symmetry correspondence as the intermediate representation, to help the pose estimation of symmetry objects. Differently, we extend 2D reflection symmetry-correspondence onto 3D point-wise paired points, to fully utilize the geometry-constraint conditions which significantly improves the performance of COPSE inference.

## 3. Method

**Problem formulation.** Our goal is to estimate the 6D pose and size of a set of unseen instances from known categories,

presented by partial point clouds. We represent the 6D object pose as a rigid-body homogeneous transformation matrix $(\boldsymbol{R}, \boldsymbol{t}) \in \boldsymbol{SE}(3)$, where rotation $\boldsymbol{R} \in \boldsymbol{SO}(3)$ and translation $\boldsymbol{t} \in \mathbb{R}^3$. $\boldsymbol{SE}(3)$ and $\boldsymbol{SO}(3)$ indicate the Lie group of 3D rigid transformations and 3D rotation, individually. And the 3D size of the object is formalized as $\boldsymbol{s} \in \mathbb{R}^3$.

We denote the original partial points of a category-known instance as $\mathcal{O}_c^{ori}(\boldsymbol{R}, \boldsymbol{t}) \in \mathbb{R}^{3 \times N}$, where $c$ is the category prior of the detected object; $N$ is the number of the valid back-projected points, and $(\boldsymbol{R}, \boldsymbol{t})$ is the object pose described in the camera frame. In order to be robust against the global scale, for all input $\mathcal{O}_c^{ori}(\boldsymbol{R}, \boldsymbol{t})$ with different dimension, we shift and scale them to the unit sphere as $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$,

$$\begin{cases} \mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n) = (\mathcal{O}_c^{ori}(\boldsymbol{R}, \boldsymbol{t}) - \boldsymbol{C}_0)/S_0 \\ \boldsymbol{C}_0 = \dfrac{1}{N} \sum_{\boldsymbol{p}_i \in \mathcal{O}_c^{ori}} \boldsymbol{p}_i \\ S_0 = \max_{\boldsymbol{p}_i \in \mathcal{O}_c^{ori}} \{\|\boldsymbol{p}_i - \boldsymbol{C}_0\|_2\} \end{cases} \tag{1}$$

where $\boldsymbol{C}_0$ is the centroid of the points $\boldsymbol{p}_i \in \mathcal{O}_c^{ori}(\boldsymbol{R}, \boldsymbol{t})$, and $S_0$ is the scale factor that transforms the original points into the unit sphere. We denote $\boldsymbol{t}_n = (t_x, t_y, t_z)$ as the object center after normalization.

**Template shape.** Concretely, our algorithm will first utilize a template shape $\mathcal{T}_c(\boldsymbol{R}_0, \boldsymbol{t}_0) \in \mathcal{M}^c$ as canonical category-specific template shape, where $\boldsymbol{R}_0$ and $\boldsymbol{t}_0$ indicate canonical orientation and translation, respectively. We denote the available categorical model collections as $\mathcal{M}^c = \{M_i^c\}_{i=1}^{N_c}$, where $N_c$ is the number of models belonging to the category $c$. The intuition of this step comes from that instances across a specific category share similarity with template shape [16]. Thus we randomly choose *a single fixed*

template shape for each category from $\mathcal{M}^c$ before the network training. Our network is trained with the template, to make it robust to data variations of categories. Additionally, the template shape $\mathcal{T}_c(\boldsymbol{R}_0, \boldsymbol{t}_0)$ is sampled by Farthest Point Sampling (FPS) algorithm [28] for preserving object geometry, into a sparse $N_k$ keypoints representation $\mathcal{K}_c(\boldsymbol{R}_0, \boldsymbol{t}_0)$.

**Symbol denotation.** We summarize the symbols here. $\mathcal{O}_c^{ori}(\boldsymbol{R}, \boldsymbol{t})$ denotes the original back-projected points from category-known instances while $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ represents the normalized one. Similarly, we use $\mathcal{K}_c(\boldsymbol{R}_0, \boldsymbol{t}_0)$ to describe the canonical template shape of the category $c$ represented by $N_k$ keypoints, and $\mathcal{K}_c(\hat{\boldsymbol{R}}, \hat{\boldsymbol{t}}_0)$ is the deformed template points predicted by 3D-OCR module. $\mathcal{O}'_c(\boldsymbol{R}, \boldsymbol{t}_n)$ indicates the symmetric points with respect to $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$. $\mathcal{G}_c^{ori}(\boldsymbol{R}, \boldsymbol{t}_n)$ is denoted the merged point set from $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ and $\mathcal{O}'_c(\boldsymbol{R}, \boldsymbol{t}_n)$. We denote $\mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ as the centered point set.

**Architecture overview.** We give an overview of our DONet framework. As in Fig. 2, the potential instance-wise pixels and the object category $c$ are predicted by using an off-the-shelf approach like Mask-RCNN [9], and the partial point cloud is then back-projected from aligned RGB-D images given camera intrinsic parameters and segmented mask. Consequently, the partial observation and the predicted object label $c$ are fed into our main network for the following prediction. Our network is composed of three modules: (1) 3D-OCR (Sec. 3.1), (2) GeoReS (Sec. 3.2), and (3) MPDE (Sec. 3.3). We first formulate a two-branch network by incorporating the 3D-OCR and GeoReS module in parallel. And the MPDE module contains voting and size regression sub-module that learn dimensional information, including centripetal offset vectors and the object size.

## 3.1. 3D Orientation-Consistent Representation

Given normalized point clouds $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$, the goal of 3D-OCR is to learn the semantic shape correspondence among intra-class instances, deforming the canonical template shape to align with the observable points.

Our network uses PointNet-like structure [30, 13] as illustrated in Fig. 2. The normalized partial observation $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ and the category-specific canonical template shape $\mathcal{K}_c(\boldsymbol{R}_0, \boldsymbol{t}_0)$ are fed into the encoders to generate shape-dependent features $\boldsymbol{F}_{\mathcal{O}}$ and $\boldsymbol{F}_{\mathcal{T}}$, respectively. And point-wise, we concatenate $\boldsymbol{F}_{\mathcal{O}}$ and $\boldsymbol{F}_{\mathcal{T}}$ with the normalized partial observation $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ to generate per-point embedding vector features, thus performing shape-guided reconstruction of $\mathcal{K}_c(\hat{\boldsymbol{R}}, \hat{\boldsymbol{t}}_0)$ under clues of geometric properties from the partial observations $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$. $\hat{\boldsymbol{R}}$ and $\hat{\boldsymbol{t}}_0$ are the estimated rotation and center, which are implicitly represented by the deformed template shape in the canonical coordinate, respectively.

Overall, this module aims at learning a parametric mapping function $f_\theta(\cdot)$ to reconstruct the 3D-OCR, namely $f_\theta(\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n), \mathcal{K}_c(\boldsymbol{R}_0, \boldsymbol{t}_0)) \rightarrow \mathcal{K}_c(\hat{\boldsymbol{R}}, \hat{\boldsymbol{t}}_0)$, where $\theta$ is optimized CNN parameters, $(\hat{\boldsymbol{R}}, \hat{\boldsymbol{t}}_0)$ are assumed as the implicit pose of the reconstructed 3D-OCR. That is to say, the reconstructed 3D-OCR $\mathcal{K}_c(\hat{\boldsymbol{R}}, \hat{\boldsymbol{t}}_0)$ implicitly characterizes the predicted orientation $\hat{\boldsymbol{R}}$ of the partial points $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$. One explanation is that the $f_\theta(\cdot)$ learns the semantic shape-correspondence between the partial observation $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ and the defined template shape $\mathcal{K}_c(\boldsymbol{R}_0, \boldsymbol{t}_0)$, to maintain semantic alignment. Then the rotation recovery $\hat{\boldsymbol{R}}$ is known as the orthogonal Procrustes problem [35] of alignment for two ordered set of keypoints $\mathcal{K}_c(\boldsymbol{R}_0, \boldsymbol{t}_0)$ and $\mathcal{K}_c(\hat{\boldsymbol{R}}, \hat{\boldsymbol{t}}_0)$.

## 3.2. Geometry-constrained Reflection Symmetry

We discuss this section under the assumption that most of the objects are symmetric or weakly symmetric. Especially, we consider two kinds of symmetry on the objects.

**Reflection symmetry.** As for reflection symmetry categories like `mug` and `laptop`, they are usually symmetric about a fixed vertical plane, treated as prior symmetry-axis constraint for easier object shape completion. Thus, given the observation $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$, there exists corresponding $\mathcal{O}'_c(\boldsymbol{R}, \boldsymbol{t}_n)$, to be symmetric with $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$.

**Rotational symmetry.** The rotational symmetric categories like `can` and `bowl` poss infinite reflection symmetry planes which hinders the network to get converged. One remedy is to rotate the partial point $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ by $180°$ around its symmetry-axis in the object frame, to generate the paired points $\mathcal{O}'_c(\boldsymbol{R}, \boldsymbol{t}_n)$. It also enables our network to reason the occluded part from the observable one, to obtain a more complete shape for subsequent dimensional estimation.

Totally, our GeoReS module learns a function $h_\gamma(\cdot)$ to enforce that the predicted mirrored points $h_\gamma(\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)) \equiv \mathcal{O}'_c(\boldsymbol{R}, \boldsymbol{t}_n)$, which aims to capture more discriminative semantic features from observation $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$; $\gamma$ is the parametric CNN parameters. We do so by using similar feature concatenation strategy as Sec. 3.1, to generate point-wise paired prediction $\hat{\mathcal{O}}'_c(\boldsymbol{R}, \boldsymbol{t}_n)$.

**Discussion of asymmetric objects.** It is noteworthy that the GeoReS module is applicable to asymmetric objects which have the global symmetric shape but asymmetric local parts. GeoReS will make an initial estimation of the shape and size of the convex hull for asymmetric objects, and the corresponding estimation errors should be further updated by the MPDE module below. Essentially, reasoning such symmetric properties of objects encourages to capture 3D semantic shape, and learn size information as discussed in Sec. 4.2.

## 3.3. Mirror-Paired Dimensional Estimation

Built upon GeoReS module, we further obtain a relatively complete object shape by merging the predicted symmetric points (orange points) $\hat{\mathcal{O}}'_c(\boldsymbol{R}, \boldsymbol{t}_n)$ and partial points (blue points) $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ as in Fig .2. Note that the predicted

symmetric points and partial points are reflected before the concatenation operation. Such a combination step generates a rough 3D object shape $\mathcal{G}_c^{ori}(\boldsymbol{R}, \boldsymbol{t}_n)$ for coarse center localization. Then, we centralize points $\mathcal{G}_c^{ori}(\boldsymbol{R}, \boldsymbol{t}_n')$, denoted as $\mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n')$,

$$
\begin{cases}
\mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n') = \mathcal{G}_c^{ori}(\boldsymbol{R}, \boldsymbol{t}_n) - \boldsymbol{C}_0' \\
\boldsymbol{C}_0' = \frac{1}{2N} \sum_{\boldsymbol{q}_i \in \mathcal{G}_c^{ori}} \boldsymbol{q}_i
\end{cases}
\tag{2}
$$

where $\boldsymbol{C}_0'$ is the computed centroid of the merged points $\boldsymbol{q}_i \in \mathcal{G}_c^{ori}(\boldsymbol{R}, \boldsymbol{t}_n)$. $\boldsymbol{t}_n'$ is the object center after centralizing. Particularly, we incorporate ground-truth symmetric points $\mathcal{O}'_c(\boldsymbol{R}, \boldsymbol{t}_n)$ and the partial observation $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ during the training phase, to prevent the unstable gradient propagation in the early phase.

**Center localization by dense voting.** Inspired by previous 2D [46, 28] and 3D [10, 29] keypoints voting methods, we treat the object center as a specific keypoint, and all the points in $\mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n')$ generate offset vectors to the potential center densely. Thus, the network optimizes the learned point-wise offset vectors, to generate the translation offset candidate set $\mathcal{D} = \{\hat{\boldsymbol{d}}_i\}_{i=1}^{2N}$, directing from the surface points $\boldsymbol{q}_i' \in \mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n')$ to the instance center $\boldsymbol{t}_n'$. Specially, we apply simple average mean to determine the final voted center $\hat{\boldsymbol{t}}_n'$ for reducing the computational expense in our practical experiments. In total, the centerd points set $\mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n')$ locates coarsely the object center, and the following point-wise dense voting searches for fine-grained potential offsets in constrained 3D space. Finally, the predicted object center $\hat{\boldsymbol{t}}$ of the original observation $\mathcal{O}_c^{ori}(\boldsymbol{R}, \boldsymbol{t})$ described in the camera frame is written as,

$$
\begin{cases}
\hat{\boldsymbol{t}} = (\hat{\boldsymbol{t}}_n' + \boldsymbol{C}_0') \cdot S_0 + \boldsymbol{C}_0 \\
\hat{\boldsymbol{t}}_n' = \frac{1}{2N} \sum_{\boldsymbol{q}_i' \in \mathcal{G}_c, \hat{\boldsymbol{d}}_i \in \mathcal{D}} (\boldsymbol{q}_i' + \hat{\boldsymbol{d}}_i)
\end{cases}
\tag{3}
$$

**Size estimation.** Obtained merging mirror-paired points $\mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n')$, the network directly regresses the normalized scale $\hat{\boldsymbol{s}}_n$. Then the actual size $\hat{\boldsymbol{s}}$ is recovered by the effect of given scale factor $S_0$, i.e. $\hat{\boldsymbol{s}} = S_0 \cdot \hat{\boldsymbol{s}}_n$. Compared to regressing only with partial points set $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$, the merging points $\mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n')$ provide a relatively complete shape prior for more accurate size estimation as in Sec. 4.2.

### 3.4. Loss Function

To efficiently train our network, we introduce and define the loss function $\mathcal{L}$ as follows,

$$
\mathcal{L} = \mathcal{L}_{rec} + \mathcal{L}_{ref} + \mathcal{L}_{cen} + \mathcal{L}_{size}
\tag{4}
$$

**Keypoint reconstruction loss.** We aim to perform the semantic shape-guided reconstruction from the observable points. Given the ground-truth 3D-OCR $\mathcal{K}_c = \{\boldsymbol{p}_i^o\}_{i=1}^{N_k}$ which consists of $N_k$ ordered keypoints $\boldsymbol{p}_i^o$ during the training phase, the network learns to reconstruct the underlying

unified representation $\hat{\mathcal{K}}_c = \{\hat{\boldsymbol{p}}_i^o\}_{i=1}^{N_k}$, the reconstruction loss is defined as:

$$
\mathcal{L}_{rec} = \frac{1}{N_k} \sum_{i=1}^{N_k} \|\boldsymbol{p}_i^o - \hat{\boldsymbol{p}}_i^o\|_1
\tag{5}
$$

Nevertheless, the $\mathcal{L}_{rec}$ encounters ambiguities in dealing with symmetric instances, which deviates the model from converging. Hence, we adopt the strategy as [43] to generate a candidate ground-truth set as follows,

$$
\mathcal{L}_{rec}^{sym} = \min_{C \in \mathcal{C}} \Big\{ \frac{1}{N_k} \sum_{i=1}^{N_k} \|C \cdot \boldsymbol{p}_i^o - \hat{\boldsymbol{p}}_i^o\|_1 \Big\}
\tag{6}
$$

where $\mathcal{C} = \{C_i\}_{i=1}^{M}$ is a set of candidate symmetry transformations. Set $\mathcal{C}$ covers $M$ equally discrete global rotational symmetries along symmetry-axis, $M = 12$ in our case.

**Symmetry prediction loss.** Our GeoReS module learns to predict potential mirrored points, that is, to localize the point-wise mirrored points based on the input points. We optimize the objective by the following loss:

$$
\mathcal{L}_{ref} = \frac{1}{N} \sum_{i=1}^{N} \|\boldsymbol{p}_i^r - \hat{\boldsymbol{p}}_i^r\|_1
\tag{7}
$$

where $N$ is the number of input points, $\boldsymbol{p}_i^r$ and $\hat{\boldsymbol{p}}_i^r$ are the ground-truth and predicted mirrored points, respectively. The network parameterizes the cues of 3D semantic shape details from the training samples.

**Centripetal offset vector loss.** The network learns the centripetal offset vector $\boldsymbol{d}_i$ pointing to the potential center from merged points $\boldsymbol{q}_i' \in \mathcal{G}_c(\boldsymbol{R}, \boldsymbol{t}_n')$, or translation offset between the surface points and object center. The learning of $\boldsymbol{d}_i$ is supervised by minimizing the predicted bias as:

$$
\mathcal{L}_{cen} = \frac{1}{2N} \sum_{i=1}^{2N} \|\boldsymbol{d_i} - \hat{\boldsymbol{d_i}}\|_1
\tag{8}
$$

where $2N$ is the number of merged points, and $\boldsymbol{d}_i$ and $\hat{\boldsymbol{d}}_i$ are the ground-truth and predicted offset vectors.

**Size loss.** In a way, the partial points $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n)$ determine the magnitude information of object center and size which is relevant to its underlying 3D shape. For better size recovery, we also take the advantage of the predicted output from the GeoReS module which constructs an explicit but simpler shape completion as discussed in Sec.3.3.

We supervise the size regression with $\ell_1$ loss:

$$
\mathcal{L}_{size} = \|\hat{\boldsymbol{s}}_n / \boldsymbol{s}_n - \boldsymbol{1}\|_1
\tag{9}
$$

where $\boldsymbol{s}_n$ represents the ground-truth normalized scale and $\hat{\boldsymbol{s}}_n$ is the predicted one of the normalized points $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}_n')$.

## 4. Experiments

**Datasets.** We conduct experiments on NOCS [43] dataset and the physical Baxter robot. Particularly, the recently introduced NOCS dataset contains six table-scale object categories including `bottle`, `bowl`, `can`, `camera`, `laptop`,

Table 1: Quantitative comparison with other competitors. We report the evaluation on mAP, accuracy on the NOCS-REAL275 datasets. And we also compare the model size of different methods. (↑): higher better, (↓): lower better.

| Dataset | Method | mAP (↑) | | | | | | Accuracy (↑) | Model size (↓) |
| | | $3D_{50}$ | $3D_{75}$ | $5°2cm$ | $5°5cm$ | $10°2cm$ | $10°5cm$ | $5°5cm$ | (M) |
|---|---|---|---|---|---|---|---|---|---|
| | NOCS [43] | 78.0 | 30.1 | 7.2 | 10.0 | 13.8 | 25.2 | - | 211.0* |
| | CASS [4] | 77.7 | - | - | 23.5 | - | 58.0 | - | - |
| REAL275 | 6-PACK [41] | - | - | - | - | - | - | 33.3 | 81.5 |
| | DeformNet [39] | 77.3 | 53.2 | 19.3 | 21.4 | 43.2 | 54.1 | 30.4 | 73.3 |
| | Ours | **80.4** | **63.7** | **24.1** | **34.8** | **45.3** | **67.4** | **49.1** | **18.5** |

Table 2: Ablation studies on different configurations for 6D object pose and size estimation on the NOCS-REAL275 dataset. (Note: FPS $N_k$ means that applied $N_k$ keypoints representation of the 3D-OCR, generated by the FPS algorithm.)

| Row | Method | | | | | | | FPS | mAP (↑) | | | | | | Accuracy (↑) |
| | 3D-OCR | GeoReS | MPDE* | MPDE | R | T | SC | | $3D_{50}$ | $3D_{75}$ | $5°2cm$ | $5°5cm$ | $10°2cm$ | $10°5cm$ | $5°5cm$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | ✓ | | | | | | | 36 | 81.1 | 55.3 | 15.9 | 23.6 | 34.9 | 59.1 | 39.4 |
| 2 | ✓ | ✓ | | | | | | 36 | **81.2** | 60.1 | 17.8 | 27.7 | 38.8 | 63.3 | 44.1 |
| 3 | ✓ | ✓ | ✓ | | | | | 36 | 80.6 | 62.6 | 20.5 | 31.7 | 39.8 | 65.3 | 46.4 |
| 4 | ✓ | ✓ | | ✓ | | | | 36 | 80.4 | **63.7** | **24.1** | **34.8** | 45.3 | 67.4 | **49.1** |
| 5 | | ✓ | | | ✓ | ✓ | | 36 | 80.6 | 62.9 | 20.8 | 29.7 | 43.6 | 64.4 | 46.3 |
| 6 | ✓ | ✓ | | | ✓ | | | 36 | 81.0 | 63.5 | 21.1 | 30.4 | 45.7 | 67.7 | 46.7 |
| 7 | ✓ | | | ✓ | | | ✓ | 36 | 80.6 | 59.5 | 19.2 | 28.9 | 41.6 | 65.2 | 45.0 |
| 8 | ✓ | ✓ | | ✓ | | | | 16 | 79.6 | 62.9 | 22.8 | 33.0 | **46.1** | **67.6** | 47.9 |
| 9 | ✓ | ✓ | | ✓ | | | | 128 | 79.5 | 59.5 | 21.5 | 32.1 | 43.7 | 66.3 | 47.5 |

and `mug`. For real-world performance evaluation, we test our method on the NOCS-REAL275 dataset which contains 2.75K real scene images in simulated real-world clutter manner. And we utilize the 3D models from ShapeNet-Core [3] for training our COPSE model. We also test our COPSE model in additional 6 different real scenes with 25 unseen instances from categories including `bowl`, `mug`, `bottle`, and `laptop`. The test scene images are not manually pose-annotated, which are just used to show qualitative results. The robotic experiments compare the real-world performance of deploying COPSE models on a real Baxter robot executing manipulation tasks. Baxter is a dual-arm collaborative robot with parallel grippers, mounted with a RealSense D435 Camera on the base. The COPSE models are trained on synthetic data of 6 categories, and tested on various real-world robotic tasks, including the tasks of grasping, pouring, and delivery. Our codes and pre-trained models will be released on GitHub.

**Evaluation metrics.** As [39], we compute the average precision of 3D Intersection-Over-Union(IoU) at 25%, 50% and 75% for object detection. And the average precision under $m°ncm$ is calculated for evaluating 6D pose recovery. Here we choose thresholds under diverse levels of difficulty at $5°2cm$, $5°5cm$, $10°2cm$, and $10°5cm$. As for the symmetric instances, only the error between the symmetry-axis is considered. And we also adopt the setting as [39, 43], treating the `mug` as symmetric objects under the specific perspective where the handle is occluded.

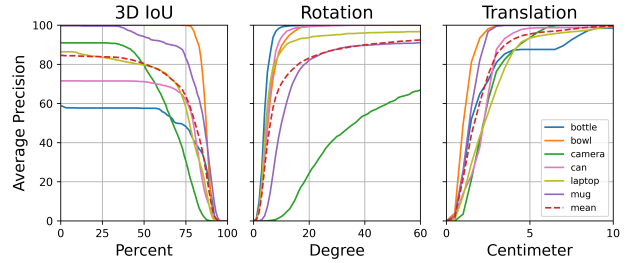**Implementation details.** Our network is exclusively



Figure 3: Average precision results of different thresholds with 3D IoU, rotation error, and translation error on the NOCS-REAL275 dataset.

trained upon synthetic depth images. To be specific, we pick 1085 object models covering six categories from ShapeNet-Core [3] and utilize the Blender software [1] to render depth images under random perspectives. Specially, we render 60 different views for each instance and back-projected partial points are disorderly sampled into 1024 points. The model is trained for 200 epochs with a batch size of 32. We initially set the learning rate as 0.0001 and decrease it by a factor of 0.75 for every 20 epoch. In robotic experiments, the estimator is implemented on the platform with an NVIDIA RTX2070 GPU and an Intel Core i7-9700K CPU, and it estimates the poses and sizes at around 10 Hz.

## 4.1. Results

**Comparisons with competitors.** We compare our approach with the RGB-D pose estimation methods [39, 4]

and the pose tracking method [41] on the NOCE-REAL275 dataset. Quantitative results are summarized in Tab. 1.

Trained upon the synthetic 3D models, all the test instances in the NOCS-REAL275 dataset are unseen ones to our network. The baseline methods require real training images, and only half of the instances in NOCS-REAL275 are unseen to their models. Even in such comparison, our approach still outperforms the state-of-the-art method [39] by a margin of 10.5% and 13.4% for the $3D_{75}$ and $5°5cm$ metric, respectively. The results also validate the efficacy of our COPSE model in sim-to-real generalization. Figure. 3 shows a more detailed analysis of 3DIoU and 6D pose error on different evaluation metrics for all 6 testing instance categories. By taking advantage of the symmetry reasoning, our approach achieves desirable rotation recovery especially for mug (reflection symmetry) which is one of the most difficult categories to estimate in previous methods [43, 39, 4]. As for the categories like bottle, bowl and can (rotational symmetry), our method easily achieves superior performance in pose recovery.

Moreover, compared to other methods as in Tab. 1, our depth-based method has a more lightweight model, which is more feasible for mobile devices. Empirically, 6-PACK [41] is a pose tracking method with higher pose accuracy versus the pose estimation methods, but it trains the model which is capable of tracking novel *category-specific* instances. We also report the mean accuracy ($<5°5cm$ percentage) of our *multi-category* model which still outperforms it by 15.8%. Although our model is trained with synthetic data, some key useful geometry in synthetic NOCS is not available, like camera, which negatively affects the performance of our COPSE model.

**Evaluation on Instance-level Task.** Our COPSE model could be easily employed in the instance-level task, by using the exact object model as the template shape. Compared with RGB(D) methods [42, 28] or depth-only method [8], our DONet still achieves good results on instance-level dataset LineMOD [11] in terms of ADD(-S) metric as below. Also, take weakly symmetric objects (e.g., cat, ape) as symmetric ones, our model still gains desirable performance, which provides evidence to support our principle of GeoReS module to handle asymmetric objects.

Table 3: Comparisons with instance-level methods of kinds of modalities on LINEMOD dataset with ADD(S) metric.

| Modality | Methods | ape | can | cat | driller | eggbox | glue |
|---|---|---|---|---|---|---|---|
| RGB | PVNet [28] | 43.6 | 95.5 | 79.3 | 96.4 | 99.1 | 95.7 |
| RGBD | DF [42] | 92 | 93 | 97 | 87 | 100 | 100 |
| synthetic D | CP+ICP [8] | 58.3 | 84.7 | 84.6 | 43.2 | 99.5 | 98.8 |
| synthetic D | Ours | 64.5 | 83.6 | 91.4 | 84.0 | 99.4 | 100 |
| synthetic D | Ours+ICP | **98.8** | **99.2** | **99.8** | **99.4** | **100** | **100** |

*It incorporates detection and pose estimation framework, not considered in our mere comparison of pose estimator.

**Qualitative results.** As for NOCS-REAL275 dataset, we qualitatively show some comparative COPSE results of ours (top row) and DeformNet [39] (bottom row), as in Fig. 4. As can be observed, our depth-based estimator that is only trained on synthetic data still achieves more accurate results than RGB-D methods. Especially, in *additional real scenes* with multi-objects, our algorithm generates accurate estimation, under which the objects are visually located within predicted bounding boxes and axis-aligned in Fig 5 (top row). Our results indicate the generalization capability of our COPSE model in real-world applications.

### 4.2. Ablation Study.

We verify the effectiveness of the key components in our model, including 3D-OCR, GeoReS, MPDE, and keypoint number. The results are summarized in Tab. 2.

**Efficacy of symmetry reasoning.** We first check the importance of symmetric properties utilization by GeoRes module. We start from a basic network, which directly outputs the 3D-OCR, voting vectors, and normalized scale based on *partial points* $\mathcal{O}_c(\boldsymbol{R}, \boldsymbol{t}'_n)$, as in **row 1**. Sequentially, the GeoReS module is added to the aforementioned basic network shown in **row 2**. The comparison results between **rows 1** and **2** illustrate that the GeoReS module is a vital part to produce overall great results. And our depth-only method without GeoReS still achieves similar performance compared to RGB-D method CASS [4], which further indicates the validity of exploiting symmetry. With symmetry reasoning, the network captures semantic-specific details for mirrored points reconstruction, from where it also explains the enhancement to the 3D-OCR module that relies on intra-class semantic shape correspondence.

**Efficacy of rough shape for object center and size estimation.** We then study the progressive improvement of object center and size estimation from roughly completed shape by the GeoReS module, compared to that of partial points. We analyze it from two aspects. First, we remove the beforehand centralizing operation in the MPDE module ("MPDE*") to merely validate the performance upon a more complete shape as in **row 3**. Comparing the results in **row 2** and **3**, estimation based on a complete shape gains overall improved performance than that of partial points. Second, we explore the importance of the provided "coarse" object center from the rough shape, by recovering the centralizing operation (full model) as in **row 4**. That is, the merged point cloud is moved to the place where its centroid is closer to the origin point, beneficial to search space constraint and localization bias reduction. The comparison results of **rows 3** and **4** indicate the necessity of the centralizing operation, yielding further improvement of 1.1% at $3D_{75}$ and 3.6% at $5°2cm$, respectively.

**3D-OCR and quaternion representation.** We investigate the usefulness of the 3D-OCR module. We replace the 3D-
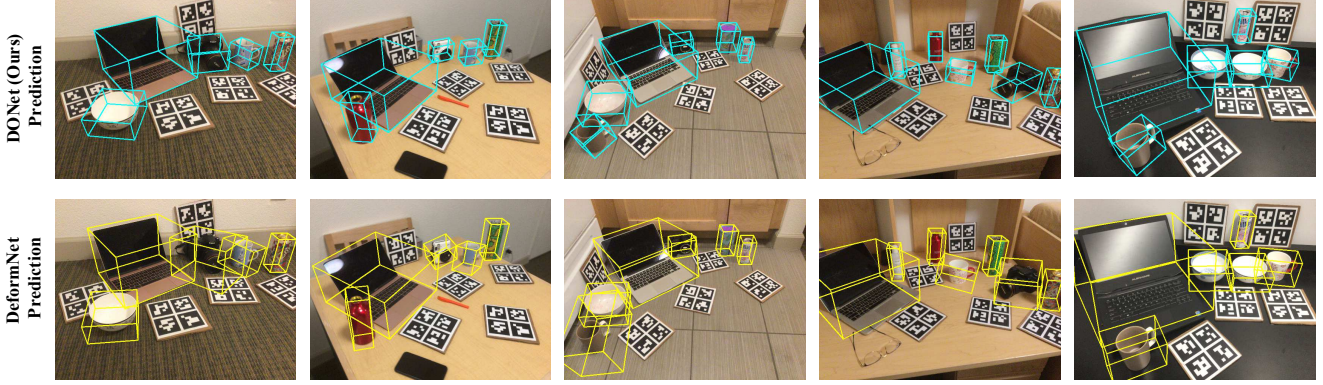
Figure 4: Qualitative comparisons between DONet and DeformNet [39]. We visualize the estimated 6D pose and size as the tight-oriented bounding box around the target instances.

OCR module from the full model with a direct quaternion regression module ("R"), as in **row 5**. In contrast to **row 4**, the performance in **row 5** decreases, which indicates that 3D-OCR which learns shape correspondence mapping is more robust and generalized than direct rotation regression. **GeoRes and object shape completion.** Also, we compare our GeoRes module and direct object completion. Concretely, we employ the same network for straightforward object completion ("SC"). All metrics in **row 7** drop consistently. This is probably because the shape completion focuses on detailed reconstruction which relies on a more complicated network, but the GeoRes module only captures the semantic shape details to ease the shape reconstruction. Poor reconstruction results from a straightforward object completion network further degrade the performance of the MPDE module.

**Voting mechanism and direct object center regression.** We replace the voting mechanism by regressing the object center $\hat{t}'_n$ ("T"). Surprisingly, direct regression has nearly the same results on $5cm$ metric, but the performance drops on a more strict $2cm$ metric, comparing **rows 4** and **6**. Therefore, the voting mechanism locates a more accurate object center.

**Verifying keypoint number.** In addition, We also explore the impact of the varied number of template points $\mathcal{K}_c(\boldsymbol{R}_0, \boldsymbol{t}_0)$ by using the full model. It is observed from the last three rows that 36 keypoints is a good trade-off for our network to learn. And the choice of 128 keypoints degrades the performance due to the larger output space, while the 16 keypoints one is too sparse to represent the geometric structure, which also results in a negative influence on the final performance.

### 4.3. Robotic Grasping Experiment

**Grasping task.** Particularly, we use 6 unseen instances from 3 classes, i.e., 2 mugs, 2 bottles, and 2 bowls. Deployed the COPSE models, Baxter takes 10 trials in grasping each object. In this experiment, our DONet is compared
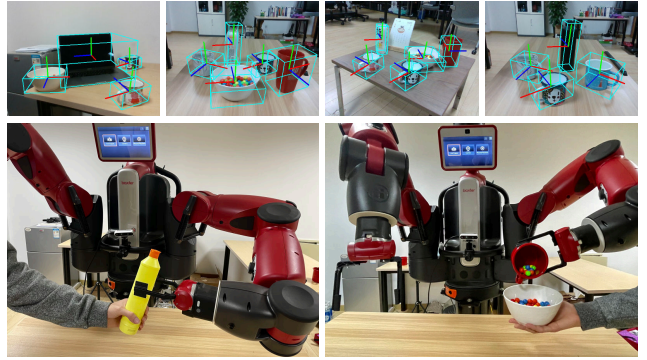


Figure 5: Here we show the estimated results given by our DONet in various real clutter environments (top row). We also perform the interaction tasks on a physical Baxter robot integrated with our DONet estimator (bottom row).

against the DeformNet [39] with the success rate of 86.7% and 68.3%, respectively. This shows the promising real-world applicability of our COPSE model.

**Delivery task.** In the delivery task, the robot interacts with the actor, trying to grasp the objects in human hands, as in Fig 5 (bottom left). We choose the testing instance bottle. The Baxter successes on 80% using our DONet in 15 trials of the delivery task, compared against that of 66.7% using DeformNet [39], validating the accuracy and real-time performance of our pose estimator.

**Pouring task.** By using our COPSE model, we conduct the pouring task, which makes the actor randomly move the bowl, while the robot follows the actor and executes pouring action at the appropriate time, as in Fig 5 (bottom right). We choose each testing instance from bowl and mug, respectively. By deploying COPSE models, Baxter makes 15 trials in pouring beans. our DONet is still compared against DeformNet [39] with the success rate of 73.3% and 53.3%, respectively. This again demonstrates the effectiveness of our COPSE model in human-robot interaction tasks.

# 5. Conclusion

We propose a novel depth-based method for category-level 6D object pose and size estimation with multiple modules, including the 3D-OCR, GeoReS, and MPDE. We also show that jointly training these modules can encourage overall performance. Finally, our method achieves state-of-the-art performance, without real-world training data. Furthermore, a physical Baxter robot integrated with our framework validates the utility in real-world robotic applications.

**Future work.** Under the inherent limitation of the depth-based method, the sensor noise and lacked discriminative details may result in ambiguities in pose recovery. Future work will consider the additional color information from RGB channels for more accurate pose and size recovery.

## References

[1] Blender software. https://www.blender.org/. 6

[2] Zhe Cao, Yaser Sheikh, and Natasha Kholgade Banerjee. Real-time scalable 6dof pose estimation for textureless objects. In *2016 IEEE International conference on Robotics and Automation (ICRA)*, pages 2441–2448. IEEE, 2016. 2

[3] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 6

[4] Dengsheng Chen, Jun Li, Zheng Wang, and Kai Xu. Learning canonical shape space for category-level 6d object pose and size estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11973–11982, 2020. 1, 2, 3, 6, 7

[5] Alvaro Collet, Dmitry Berenson, Siddhartha S Srinivasa, and Dave Ferguson. Object recognition and full pose registration from a single image for robotic manipulation. In *2009 IEEE International Conference on Robotics and Automation*, pages 48–55. IEEE, 2009. 1

[6] Xinke Deng, Yu Xiang, Arsalan Mousavian, Clemens Eppner, Timothy Bretl, and Dieter Fox. Self-supervised 6d object pose estimation for robot manipulation. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3665–3671. IEEE, 2020. 1

[7] Christopher Funk, Seungkyu Lee, Martin R Oswald, Stavros Tsogkas, Wei Shen, Andrea Cohen, Sven Dickinson, and Yanxi Liu. 2017 iccv challenge: Detecting symmetry in the wild. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1692–1701, 2017. 3

[8] Ge Gao, Mikko Lauri, Yulong Wang, Xiaolin Hu, Jianwei Zhang, and Simone Frintrop. 6d object pose regression via supervised learning on point clouds. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3643–3649. IEEE, 2020. 7

[9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 4

[10] Yisheng He, Wei Sun, Haibin Huang, Jianran Liu, Haoqiang Fan, and Jian Sun. Pvn3d: A deep point-wise 3d keypoints voting network for 6dof pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11632–11641, 2020. 2, 5

[11] Stefan Hinterstoisser, Vincent Lepetit, Slobodan Ilic, Stefan Holzer, Gary Bradski, Kurt Konolige, and Nassir Navab. Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes. In *Asian conference on computer vision*, pages 548–562. Springer, 2012. 2, 7

[12] Yinlin Hu, Pascal Fua, Wei Wang, and Mathieu Salzmann. Single-stage 6d object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2930–2939, 2020. 1

[13] Zitian Huang, Yikuan Yu, Jiawen Xu, Feng Ni, and Xinyi Le. Pf-net: Point fractal network for 3d point cloud completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7662–7670, 2020. 4

[14] Jarmo Ilonen, Jeannette Bohg, and Ville Kyrki. Three-dimensional object reconstruction of symmetric objects by fusing visual and tactile sensing. *The International Journal of Robotics Research*, 33(2):321–341, 2014. 3

[15] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946, 2015. 2

[16] Nilesh Kulkarni, Abhinav Gupta, and Shubham Tulsiani. Canonical surface mapping via geometric cycle consistency. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2202–2211, 2019. 3

[17] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. Epnp: An accurate o (n) solution to the pnp problem. *International journal of computer vision*, 81(2):155, 2009. 2

[18] Yan Li, Leon Gu, and Takeo Kanade. Robustly aligning a shape model and its application to car alignment of unknown pose. *IEEE transactions on pattern analysis and machine intelligence*, 33(9):1860–1876, 2011. 2

[19] Yi Li, Gu Wang, Xiangyang Ji, Yu Xiang, and Dieter Fox. Deepim: Deep iterative matching for 6d pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 683–698, 2018. 1

[20] Zhigang Li, Gu Wang, and Xiangyang Ji. Cdpn: Coordinates-based disentangled pose network for real-time rgb-based 6-dof object pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7678–7687, 2019. 2

[21] David G Lowe. Object recognition from local scale-invariant features. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1150–1157. Ieee, 1999. 2

[22] Lucas Manuelli, Wei Gao, Peter Florence, and Russ Tedrake. kpam: Keypoint affordances for category-level robotic manipulation. *arXiv preprint arXiv:1903.06684*, 2019. 1

[23] Eric Marchand, Hideaki Uchiyama, and Fabien Spindler. Pose estimation for augmented reality: a hands-on survey. *IEEE transactions on visualization and computer graphics*, 22(12):2633–2651, 2015. 1

[24] Rajendra Nagar and Shanmuganathan Raman. Reflection symmetry axes detection using multiple model fitting. *IEEE Signal Processing Letters*, 24(10):1438–1442, 2017. 3

[25] Markus Oberweger, Mahdi Rad, and Vincent Lepetit. Making deep heatmaps robust to partial occlusions for 3d object pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 119–134, 2018. 2

[26] Kiru Park, Timothy Patten, and Markus Vincze. Pix2pose: Pixel-wise coordinate regression of objects for 6d pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7668–7677, 2019. 1

[27] Georgios Pavlakos, Xiaowei Zhou, Aaron Chan, Konstantinos G Derpanis, and Kostas Daniilidis. 6-dof object pose from semantic keypoints. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 2011–2018. IEEE, 2017. 2

[28] Sida Peng, Yuan Liu, Qixing Huang, Xiaowei Zhou, and Hujun Bao. Pvnet: Pixel-wise voting network for 6dof pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4561–4570, 2019. 1, 2, 4, 5, 7

[29] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9277–9286, 2019. 5

[30] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 4

[31] Mahdi Rad and Vincent Lepetit. Bb8: A scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3828–3836, 2017. 2

[32] Caner Sahin, Guillermo Garcia-Hernando, Juil Sock, and Tae-Kyun Kim. Instance-and category-level 6d object pose estimation. In *RGB-D Image Analysis and Processing*, pages 243–265. Springer, 2019. 1

[33] Caner Sahin, Guillermo Garcia-Hernando, Juil Sock, and Tae-Kyun Kim. A review on object pose recovery: from 3d bounding box detectors to full 6d pose estimators. *Image and Vision Computing*, page 103898, 2020. 1

[34] Caner Sahin and Tae-Kyun Kim. Category-level 6d object pose recovery in depth images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018. 1, 3

[35] Peter H Schönemann. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):1–10, 1966. 4

[36] Chen Song, Jiaru Song, and Qixing Huang. Hybridpose: 6d object pose estimation under hybrid representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 431–440, 2020. 3

[37] Zhiqiang Sui, Zheming Zhou, Zhen Zeng, and Odest Chadwicke Jenkins. Sum: Sequential scene understanding and manipulation. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3281–3288. IEEE, 2017. 1

[38] Bugra Tekin, Sudipta N Sinha, and Pascal Fua. Real-time seamless single shot 6d object pose prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 292–301, 2018. 2

[39] Meng Tian, Marcelo H Ang Jr, and Gim Hee Lee. Shape prior deformation for categorical 6d object pose and size estimation. *arXiv preprint arXiv:2007.08454*, 2020. 1, 2, 3, 6, 7, 8

[40] Jonathan Tremblay, Thang To, Balakumar Sundaralingam, Yu Xiang, Dieter Fox, and Stan Birchfield. Deep object pose estimation for semantic robotic grasping of household objects. *arXiv preprint arXiv:1809.10790*, 2018. 1

[41] Chen Wang, Roberto Martín-Martín, Danfei Xu, Jun Lv, Cewu Lu, Li Fei-Fei, Silvio Savarese, and Yuke Zhu. 6-pack: Category-level 6d pose tracker with anchor-based keypoints. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 10059–10066. IEEE, 2020. 3, 6, 7

[42] Chen Wang, Danfei Xu, Yuke Zhu, Roberto Martín-Martín, Cewu Lu, Li Fei-Fei, and Silvio Savarese. Densefusion: 6d object pose estimation by iterative dense fusion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3343–3352, 2019. 1, 2, 7

[43] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J Guibas. Normalized object coordinate space for category-level 6d object pose and size estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2642–2651, 2019. 1, 2, 3, 5, 6, 7

[44] Zhaozhong Wang, Zesheng Tang, and Xiao Zhang. Reflection symmetry detection using locally affine invariant edge correspondence. *IEEE Transactions on Image Processing*, 24(4):1297–1301, 2015. 3

[45] Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Unsupervised learning of probably symmetric deformable 3d objects from images in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1–10, 2020. 3

[46] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. *arXiv preprint arXiv:1711.00199*, 2017. 2, 5

[47] Tianfan Xue, Jianzhuang Liu, and Xiaoou Tang. Symmetric piecewise planar object reconstruction from a single image. In *CVPR 2011*, pages 2577–2584. IEEE, 2011. 3

[48] Sergey Zakharov, Ivan Shugurov, and Slobodan Ilic. Dpod: 6d pose object detector and refiner. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1941–1950, 2019. 2