

## CAPTRA: Category-level Pose Tracking for Rigid and Articulated Objects from Point Clouds

Yijia Weng<sup>1\*</sup> He Wang<sup>1,2,5\*†</sup> Qiang Zhou<sup>4</sup> Yuzhe Qin<sup>3</sup> Yueqi Duan<sup>2</sup>  
 Qingnan Fan<sup>2,6</sup> Baoquan Chen<sup>1</sup> Hao Su<sup>3</sup> Leonidas J. Guibas<sup>2</sup>  
<sup>1</sup>CFCS, Peking University <sup>2</sup>Stanford University <sup>3</sup>UCSD  
<sup>4</sup>Shandong University <sup>5</sup>Beijing Institute for General AI <sup>6</sup>Tencent AI Lab

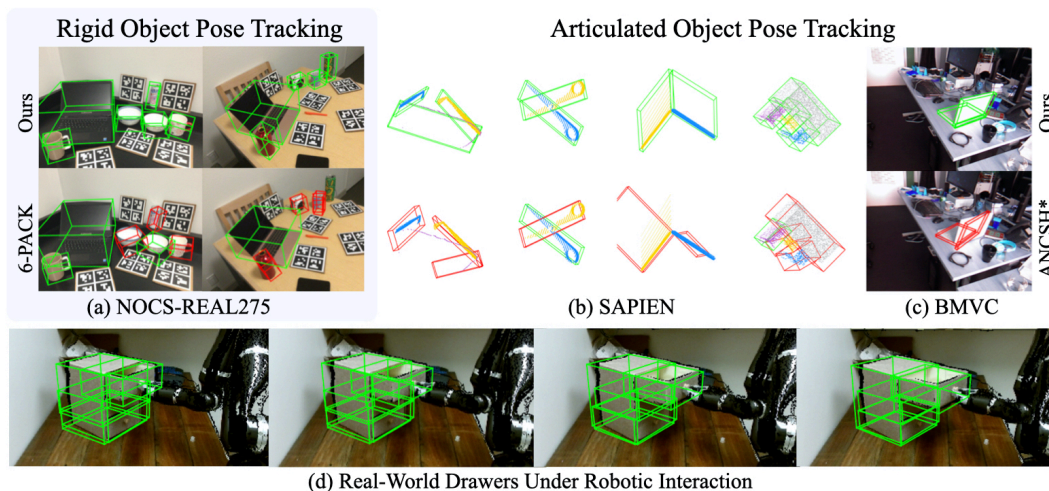


Figure 1. Our method tracks 9DoF category-level poses (3D rotation, 3D translation, and 3D size) of novel rigid objects as well as parts in articulated objects from live point cloud streams. We demonstrate: (a) our method can reliably track rigid object poses from the challenging NOCS-REAL275 dataset [29]; (b) our method can perfectly track articulated objects with big global and articulated motions from the SAPIEN datasets[34]; (c)(d) trained only on SAPIEN, our model can directly generalize to novel real laptops from BMVC dataset[18], and novel real drawers under robotic interaction. In all cases, our method significantly outperforms the previous state-of-the-arts and baselines. Here we visualize the estimated 9DoF poses as 3D bounding boxes: green boxes indicate in tracking whereas red boxes indicate off tracking.

### Abstract

In this work, we tackle the problem of category-level on-line pose tracking of objects from point cloud sequences. For the first time, we propose a unified framework that can handle 9DoF pose tracking for novel rigid object instances as well as per-part pose tracking for articulated objects from known categories. Here the 9DoF pose, comprising 6D pose and 3D size, is equivalent to a 3D amodal bounding box representation with free 6D pose. Given the depth point cloud at the current frame and the estimated pose from the last frame, our novel end-to-end pipeline learns to accurately update the pose. Our pipeline is composed of three modules: 1) a pose canonicalization module that normal-

izes the pose of the input depth point cloud; 2) RotationNet, a module that directly regresses small interframe delta rotations; and 3) CoordinateNet, a module that predicts the normalized coordinates and segmentation, enabling analytical computation of the 3D size and translation. Leveraging the small pose regime in the pose-canonicalized point clouds, our method integrates the best of both worlds by combining dense coordinate prediction and direct rotation regression, thus yielding an end-to-end differentiable pipeline optimized for 9DoF pose accuracy (without using non-differentiable RANSAC). Our extensive experiments demonstrate that our method achieves new state-of-the-art performance on category-level rigid object pose (NOCS-REAL275 [29]) and articulated object pose benchmarks (SAPIEN [34], BMVC [18]) at the fastest FPS  $\sim 12$ .

\*: equal contributions, †: corresponding author  
 Project page: <https://yijiaiweng.github.io/CAPTRA>

## 1. Introduction

Object pose estimation is crucial for a variety of computer vision and robotics applications, such as 3D scene understanding, robotic manipulation and augmented reality. The majority of object pose estimation works, *e.g.*, [35, 24], mainly lie in instance-level estimation, where the task is to estimate poses for objects from a small set of a priori known instances, thus preventing them from perceiving the poses of the vast diversity of objects in our daily life. To mitigate this limitation, Wang *et al.* [29] proposed to generalize the instance-level 6DoF (Degree of Freedom) object pose estimation problem to a category-level 9DoF pose estimation problem that takes into account the traditional 6D object pose (rotation, translation) as well as 3D object size. The proposed method in [29] can handle novel object instances in known categories without requiring CAD models of the objects. Going beyond rigid objects and in the same spirit, Li *et al.* [14] proposed to estimate category-level per-part 9DoF poses for articulated objects, such as laptops, drawers and eyeglasses.

While most of the existing category-level pose estimation works focus on single-frame estimation, we believe that temporally smooth pose tracking is more useful for many robotics applications, *e.g.*, instant feedback control, as well as AR applications. In this work, we tackle a problem named CAPTRA — *CA*tegorY-level *PO*se *TR*acking for *R*igid and *A*rticulated *O*bjects, from a live point cloud stream. Given an initial object pose at the first frame, our task is to continuously track the 9DoF pose for rigid objects or each individual rigid part of an articulated object. The most related work to ours is 6-PACK [28], which tackles the problem of category-level 6D pose tracking only for rigid objects (see the related work section for detailed comparisons).

To accurately track 9DoF poses, we consider two types of approaches: coordinate-based approaches widely used in object pose [2, 29, 14] and camera pose estimation [4] and direct pose regression as in [35, 32]. These two approaches both have pros and cons. Coordinate-based methods, which predict dense object coordinates followed by a RANSAC-based pose fitting, are generally more accurate and robust, especially on rotation estimation [27], benefiting from outlier removal in RANSAC. However, RANSAC-based pose fitting is non-differentiable and time-consuming, which often leads to a bottleneck in its running speed. In contrast, direct pose regression performs an end-to-end pose prediction, thus can achieve very high running speed, at the cost of being more error-prone.

In this work, we seek to take the best of both worlds and build **an end-to-end differentiable pipeline for accurate and fast pose tracking**. To enable highly accurate pose estimation, we propose to jointly canonicalize the input and output spaces of this estimation problem by transforming the point clouds using the inverse poses from the previous

frame. The produced **pose-canonicalized point clouds** feature near identical poses regarding the object/part, whose poses are more regression-friendly. We thus propose **RotationNet**, a PointNet++ [23] based neural network, that directly regresses the small remained rotations. Due to the ambiguity between occlusion and center translation in the partial depth observations, we found scale and translation regression still challenging. We instead propose to build **CoordinateNet** to predict dense normalized coordinates, which contain more accurate information about translation and object size due to their awareness of the category-level shape prior. Combining the outputs from both networks, we can analytically compute sizes and translations, yielding an end-to-end differentiable pipeline optimized for 9DoF pose accuracy without using non-differentiable RANSAC.

By harnessing both approaches, our proposed method gains significant performance improvement on the category-level rigid object pose benchmark and articulated object pose benchmarks. On the NOCS-REAL275 dataset [29], we outperform 6-PACK [28], the previous state-of-the-art, by 40.03% absolute improvement on the mean accuracy of  $5^{\circ}5\text{cm}$  and 10.52% absolute improvement on the mean IoU metric. On the SAPIEN articulated object dataset [34], we are the first to perform tracking and outperform the single-frame articulated pose estimation baseline, constructed using ANCSH [14] and ground truth segmentation masks, by a large margin, *e.g.*, around 20 points on mean accuracy  $5^{\circ}5\text{cm}$  in the challenging eyeglasses category. On novel real laptop trajectories from the BMVC dataset [18], we achieve the best performance directly generalized from SAPIEN without further fine-tuning. Finally, our extensive experiments further demonstrate the robustness of our tracking method to pose errors and achieve the fastest speed ( $\sim 12$  FPS) among all methods.

## 2. Related Works

**Category-Level Object Pose Estimation** To define category-level poses of novel object instances, Wang *et al.* [29] proposed Normalized Object Coordinate Space (NOCS) as a category-specific canonical reference frame for rigid objects. The objects from the same category in NOCS are consistently aligned to a category-level canonical orientation. These objects are further zero-centered and uniformly scaled so that their tight bounding boxes are centered at the origin of NOCS with a diagonal length of 1. Li *et al.* [14] extended the definition of NOCS to rigid parts in articulated objects and proposed Normalized Part Coordinate Space (NPCS), which is a part-level canonical reference frame (see appendix A for a detailed introduction). Several works have been improving [29] via leveraging analysis-by-synthesis and shape generative models as in [7, 6] and learnable deformation as in [26]. Most of these methods leverage RANSAC for pose fitting, which prohibits their pipelines from being end-to-end differentiable,

potentially rendering those methods sub-optimal. Although several works have proposed differentiable RANSAC layers to bridge this gap, *e.g.*, DSAC [2], DSAC++ [3], we will show that our method performs better than these methods without using RANSAC.

**Category-Level Object Pose Tracking** As the only existing work in this field, Wang *et al.* [28] proposed a 6D Pose Anchor-based Category-level Keypoint tracker (6-PACK) by tracking keypoints in RGB-D videos. 6-PACK first employs an attention mechanism over anchors and then generates keypoints in an unsupervised manner, which are used to compute interframe pose changes. It is important to note several key differences between 6-PACK and our work: 1) 6-PACK is designed only for rigid objects and cannot handle articulated objects; 2) 6-PACK only estimates the 6D pose containing rotation and translation and omits the important 3D size estimation required to obtain the 3D amodal object bounding boxes.

As special cases of category-level articulated object pose tracking, hand and human pose tracking problems are very popular due to their broad applications [21, 31, 20, 12, 11, 36, 1]. However, the developed methods leverage domain-specific knowledge of hand and human body, thus prevent them from being applied to generic articulated objects.

**Instance-Level 6D Pose Tracking** Instance-level pose tracking works track the poses of known object instances. Classic methods, *e.g.*, ICP-based tracking [38], explicitly rely on the object CAD models. Some recent works [8, 9, 33, 13, 9] use particle filtering to estimate the posterior of object poses. Other methods measure the discrepancy between the current observation and the previous state, and perform tracking via optimization [25, 22]. The most relevant works to ours are delta pose based methods [15, 32], which perform tracking by regressing the pose change between consecutive frames using deep neural networks.

### 3. Problem Definition and Notations

In this paper, we target at the problem of tracking the 9DoF poses of rigid or articulated objects from known categories. We follow the category-level rigid object and part pose definition in [29, 14] and adopt the assumption in [14] that the number of rigid parts  $M$  is known and constant for all the objects in a known category, where  $M > 1$  indicates an articulated object category, and  $M = 1$  indicates a rigid object category. Without loss of generality, we only describe the notations for articulated object pose tracking. For a point cloud  $X = \{x_i \in \mathbb{R}^3\}_{i=1}^N$  containing object instance  $O = \{C^{(j)}\}_{j=1}^M$ , where  $N$  is the number of points and  $C^{(j)} \subset X$  represents points of the  $j$ -th rigid part, we denote category-level part pose as  $\mathcal{P}^{(j)} = \{d^{(j)}, R^{(j)}, T^{(j)}\}$ , where  $d^{(j)} \in \mathbb{R}^3$  is 3D size,  $R^{(j)} \in SO(3)$  is rotation, and  $T^{(j)} \in \mathbb{R}^3$  is translation.

Our problem is then defined as follows: Given a live

stream of depth point clouds  $\{X_t\}_{t \geq 0}$  containing object instance  $O$  along with its per-part pose initialization  $\{\mathcal{P}_0^{(j)}\}_j$ , our objective is to track its part poses  $\{\mathcal{P}_t^{(j)}\}_j$  in an online manner, where we process the point clouds and estimate the poses for all the frames  $t > 0$ . In other words, at frame  $t+1$ , given the estimated  $\{\mathcal{P}_t^{(j)}\}_j$  from frame  $t$  and the depth point cloud  $X_{t+1}$ , our system needs to estimate  $\{\mathcal{P}_{t+1}^{(j)}\}_j$ .

## 4. End-to-end Differentiable Pose Tracking

In this section, we introduce our approach in detail. We present the pose canonicalization module in Section 4.1, and pose tracking in Section 4.2, which includes the proposed RotationNet module and CoordinateNet module, finally, we describe our training protocol in Section 4.3. The entire framework is differentiable and end-to-end trained, without any pre- or post- processing.

### 4.1. Pose Canonicalization

Inspired by [29], we factorize the 9DoF pose  $\mathcal{P}^{(j)}$  prediction into a 7DoF similarity transformation  $\mathcal{T}_t^{(j)} \in \text{Sim}(3)$  estimation problem and a 3D aspect ratio  $\hat{d}^{(j)}$  estimation problem. Formally, we define the per-part 1D uniform scale as  $s^{(j)} = \|d^{(j)}\|$  and 3D aspect ratio as  $\hat{d}^{(j)} = d^{(j)}/s^{(j)}$  so that  $d^{(j)} = s^{(j)}\hat{d}^{(j)}$ . We can then represent  $\mathcal{T}^{(j)} = \{s^{(j)}, R^{(j)}, T^{(j)}\}$ .

To estimate  $\mathcal{T}_{t+1}^{(j)}$  from the observed point cloud  $X_{t+1}$ , there are two types of approaches. One way is to train a neural network for direct pose regression. Another way is to estimate the normalized coordinates  $Y_{t+1}^{(j)}$  of  $C_{t+1}^{(j)}$ , which satisfy  $C_{t+1}^{(j)} = s_{t+1}^{(j)}R_{t+1}^{(j)}Y_{t+1}^{(j)} + T_{t+1}^{(j)}$ , and then compute the  $\mathcal{T}_{t+1}^{(j)}$  using the Umeyama algorithm [27] along with RANSAC, thus the 3D aspect ratios  $\hat{d}^{(j)}$  can be estimated using the axis range ( $|x|_{max}, |y|_{max}, |z|_{max}$ ) of  $Y_{t+1}^{(j)}$ .

In our framework, to simplify the learning problem of mapping the input  $X_{t+1}$  to the output  $\mathcal{T}_{t+1}^{(j)}$ , we propose to canonicalize both its input and output spaces using  $\mathcal{T}_t^{(j)}$ , which allows to further combine the two aforementioned methods.

**Definition** (Pose-canonicalized point cloud). *Pose-canonicalized point cloud  $Z_{t+1}^{(j)}$  with respect to part  $j$  in the input point cloud  $X_{t+1}$  is defined as the product of the inverse transformation of  $\mathcal{T}_t^{(j)}$  and  $X_{t+1}$ , namely  $Z_{t+1}^{(j)} = \left(R_t^{(j)}\right)^{-1} \left(X_{t+1} - T_t^{(j)}\right) / s_t^{(j)}$ .*

We observe that for the learning problem that maps  $X_{t+1}$  to  $\mathcal{T}_{t+1}^{(j)}$ , by canonicalizing its input  $X_{t+1}$  to pose-canonicalized point cloud  $Z_{t+1}^{(j)}$ , its output  $\mathcal{T}_{t+1}^{(j)}$  will be canonicalized to  $\hat{\mathcal{T}}_{t+1}^{(j)} = \{\hat{s}_{t+1}^{(j)}, \hat{R}_{t+1}^{(j)}, \hat{T}_{t+1}^{(j)}\}$  correspondingly, where  $\hat{s}_{t+1}^{(j)} \approx 1$ ,  $\hat{R}_{t+1}^{(j)} \approx I$ ,  $\hat{T}_{t+1}^{(j)} \approx 0$ . (See appendix B for the proof.)

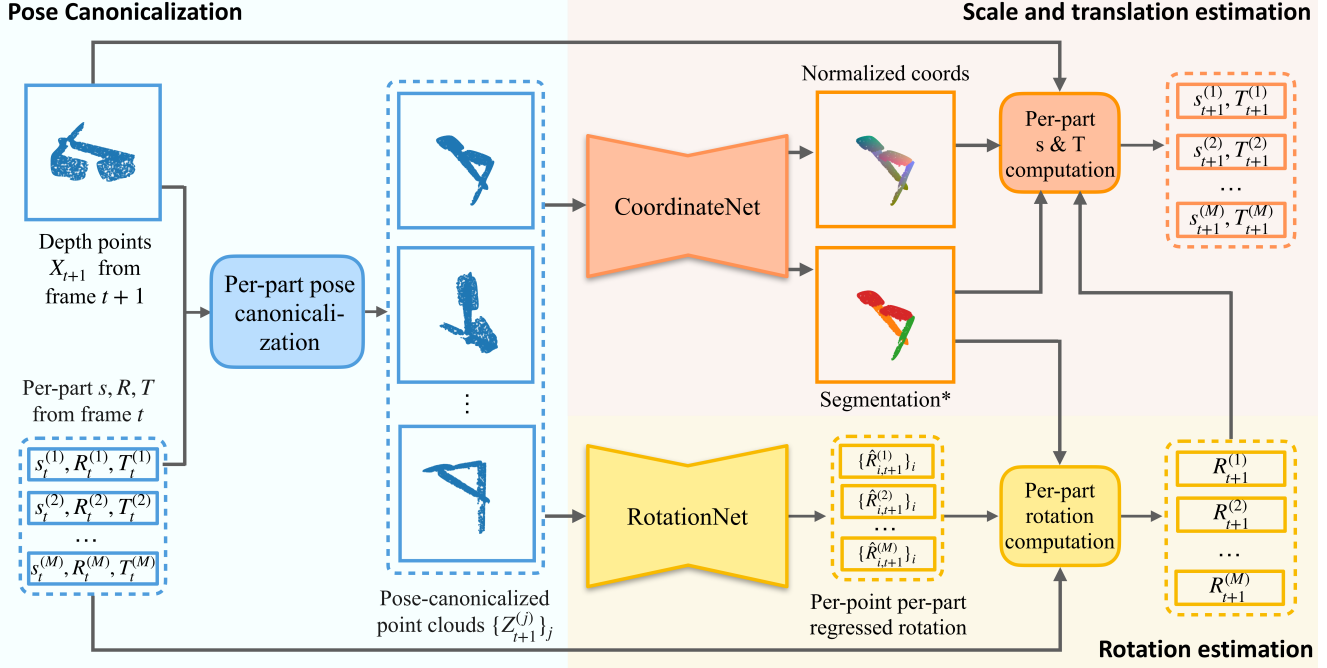


Figure 2. **Our end-to-end differentiable pose tracking pipeline** takes as inputs a depth point cloud of an  $M$ -part object along with its per-part scales, rotations, and translations estimated from the last frame. We first adopt per-part pose canonicalization to transform the depth points using the inverse estimated pose and generate  $M$  pose canonicalized point clouds. The canonicalized point clouds will be fed into RotationNet for per-part rotation estimation, as well as CoordinateNet for part segmentation and normalized coordinate predictions, which are used to compute the updated scales and translations. When RGB images are available, segmentation can be replaced by the results from the off-the-shelf image detectors for better accuracy. Such a pipeline can be naturally adapted to rigid objects when  $M = 1$ .

Note that  $\mathcal{T}_{t+1}^{(j)}$  can be expressed using  $\hat{\mathcal{T}}_{t+1}^{(j)}$  and  $\mathcal{T}_t^{(j)}$ , namely,  $s_{t+1} = s_t^{(j)} \hat{s}_t^{(j)}$ ,  $R_{t+1}^{(j)} = R_t^{(j)} \hat{R}_t^{(j)}$ ,  $T_{t+1}^{(j)} = s_t^{(j)} R_t^{(j)} \hat{T}_t^{(j)} + T_t^{(j)}$ . Now the pose estimation problem has been transformed and canonicalized to learning a mapping from  $Z_{t+1}^{(j)}$  to  $\hat{\mathcal{T}}_{t+1}^{(j)}$ . The input  $Z_{t+1}^{(j)}$  contains  $\hat{C}_{t+1}^{(j)} = (R_t^{(j)})^{-1} (C_{t+1}^{(j)} - T_t^{(j)}) / s_t^{(j)}$  that aligns well across different frames and the output space is quite constrained around an identity transformation. In this way, we simultaneously canonicalize the input point cloud space and the output pose space. By doing so, we significantly simplify the regression task, yielding much improved pose estimation accuracy and better generalizability across different instances. Note that turning the estimation of  $\mathcal{T}_{t+1}^{(j)}$  into  $\hat{\mathcal{T}}_{t+1}^{(j)}$  is closely related to estimating the interframe delta 6D pose widely used in instance-level 6D pose estimation and tracking works [15, 32]. To be more specific, we estimate interframe 7D delta transformations in a category-level canonical frame, *i.e.*, NOCS for rigid objects and NPCS for parts, whereas none of the input and output spaces of delta pose estimation in [15, 32] is canonical.

#### 4.2. Pose Tracking

Taking the pose-canonicalized point cloud  $Z_{t+1}^{(j)}$  as input, we learn a RotationNet (see Section 4.2.1) that directly regresses  $\hat{R}_{t+1}^{(j)}$  and then recovers  $R_{t+1}^{(j)} = R_t^{(j)} \hat{R}_{t+1}^{(j)}$ . Since

$\hat{R}_{t+1}^{(j)}$  lies in a small neighborhood around  $I_{3 \times 3}$ , the regression can be done with high accuracy. However, we experimentally find directly regressing  $\hat{s}_t^{(j)}$  and  $\hat{T}_t^{(j)}$  still difficult, due to the translation ambiguity caused by incompleteness of the partial observations  $Z_{t+1}^{(j)}$ . Imagine a pencil with one end occluded, the length of the pencil cannot be determined, thus making prediction of its center translation highly ambiguous. Although certain cues, *e.g.*, object symmetry, may help relieve the ambiguity, our experiments show that given partial observations of asymmetric objects, regressing  $\hat{T}_t^{(j)}$  still remains challenging. In contrast, our experiments reveal that predicting normalized coordinates  $Y_{t+1}^{(j)}$  from  $Z_{t+1}^{(j)}$  is quite successful on all rigid and articulated objects, which implicitly estimates  $\hat{s}_t^{(j)}$  and  $\hat{T}_t^{(j)}$ . The reason for this success is that the normalized coordinates  $Y_{t+1}^{(j)}$  capture the category-wise prior and enforces a zero-centered frame, thus making the translation estimation more well-considered and accurate than direct regression. We therefore devise a CoordinateNet to segment  $C_{t+1}^{(j)}$  from  $X_{t+1}^{(j)}$  and predict  $Y_{t+1}^{(j)}$  (see Section 4.2.2).

By combining the RotationNet and CoordinateNet’s outputs and knowing  $C_{t+1}^{(j)} = s_{t+1}^{(j)} R_{t+1}^{(j)} Y_{t+1}^{(j)} + T_{t+1}^{(j)}$ , we can analytically compute  $s_{t+1}^{(j)}$  and  $T_{t+1}^{(j)}$  via the Umeyama algorithm [27] (assume  $R$  is given). Usually a non-differentiable



RANSAC is needed when using Umeyama algorithm as in [29, 14] due to the multi-modal noises in the predicted  $Y_{t+1}^{(j)}$ . Thanks to the pose canonicalization, we find that our  $Y_{t+1}^{(j)}$  predictions are very successful and RANSAC only brings limited improvement to our predictions (see Sec. 5.6).

Being free from the non-differentiable RANSAC step, our end-to-end differentiable pose tracking pipeline can be straightly optimized for pose accuracy, enforce pose losses (e.g., IoU loss) directly at its outputs (see Section 4.3), and improve its running speed.

#### 4.2.1 Rotation Estimation

**RotationNet** To predict  $\{\hat{R}_{t+1}^{(j)}\}_j$  for each individual part, we devise a point cloud based deep neural network, RotationNet, that takes as inputs the pose-canonicalized point clouds  $\{Z_{t+1}^{(j)}\}_j$  with respect to each individual part  $j$ . Built upon PointNet++ [23], RotationNet regresses per-point per-part rotations  $\{\hat{R}_{i,t}^{(j)}\}_{i,j}$  in the form of the 6D continuous rotation representation [39]. Note that the PointNet++ is not deterministic since it uses random further point sampling in both set abstraction and ball query operations, thus resulting in difficulties achieving convergence on accurate regression tasks. To suppress noise, we average across the rotation predictions using the Euclidean mean [19] from points on part  $j$  to obtain the final prediction  $\hat{R}_{t+1}^{(j)}$ .

For symmetric objects such as bowls, rotation ambiguity exists around their symmetric axis. See appendix C.2 for how we supervise rotation for them.

**Training and Inference** At training time, we enforce a per-point mean square loss for points inside the ground truth mask  $m_{t+1}^{(j)}$ . At test time, the mask comes from the predicted part segmentation from CoordinateNet.

#### 4.2.2 Scale and Translation Estimation

**CoordinateNet** To estimate  $\{Y_{t+1}^{(j)}\}_j$ , we devise CoordinateNet that takes as input the pose-canonicalized point cloud  $Z_{t+1}^{(1)}$  with respect to the first part ( $j = 1$ ) and predicts its per-point part segmentation and per-point per-part normalized coordinates  $\{Y_{i,t+1}^{(j)}\}_{i,j}$ . Note that pose-canonicalized point clouds with respect to different parts share the same segmentation and normalized coordinates; thus, we only need to take  $Z_{t+1}^{(1)}$  as the input.

Built upon a PointNet++ segmentation network, CoordinateNet branches into two heads after the final feature propagation layers: one head for segmentation and the other for normalized coordinate prediction. For the segmentation head, we use relaxed IoU loss [37]. For the normalized coordinate head, we predict class-aware normalized coordinates, similar to [29, 14]. During training, we enforce an RMSE loss on the points inside the ground truth part masks. At test time, we use predicted masks to select coordinate predictions from  $M$  parts.

For symmetric objects, see appendix C.3 for how we handle ambiguity in their normalized coordinates.

**Per-part Scale and Translation Computing** Knowing the dense correspondence between  $Y_{t+1}^{(j)}$  and  $C_{t+1}^{(j)}$  and assuming  $R_{t+1}^{(j)}$  is given by RotationNet, we can analytically compute  $s_{t+1}^{(j)}$  and  $T_{t+1}^{(j)}$  via the Umeyama algorithm [27]. See appendix C.4 for further details.

### 4.3. Training Protocol

**Training Data Generation** To train CoordinateNet and RotationNet, we need paired data of pose-canonicalized point clouds and their corresponding ground truth poses. We propose to generate the training data on-the-fly without using any real temporal data. For a depth point cloud  $X$  and part  $j$  in it, we perturb its per-part ground truth scale, rotation, and translation by adding random Gaussian noise to them, namely  $s'^{(j)} = s^{(j)}(1 + n_s)$ ,  $R'^{(j)} = R^{(j)}R_{\text{rand}}$ ,  $T'^{(j)} = T^{(j)} + n_T$ , where  $n_s \sim \mathcal{N}(0, \sigma_s)$ ,  $R_{\text{rand}}$  is a rotation matrix with a random axis and a random angle  $n_\theta \sim \mathcal{N}(0, \sigma_\theta)$ , and  $n_T$  is a 3D vector with a random direction and a random length  $t \sim \mathcal{N}(0, \sigma_T)$ . We then generate the pose-canonicalized point clouds and compute their corresponding ground truth.

**Pose Losses for RotationNet and CoordinateNet** For  $s_{t+1}^{(j)}$ ,  $R_{t+1}^{(j)}$ ,  $T_{t+1}^{(j)}$ , their predictions are end-to-end differentiable; we thus propose to additionally enforce pose losses directly on these predictions. We use RMSE loss for supervising scale error  $L_{\text{scale}}$  and translation error  $L_{\text{trans}}$ . To directly improve the final 3D IoU predictions, we incorporate a corner loss  $L_{\text{corner}}$  [16], defined as the corresponding per-vertex distances between the ground truth bounding box in the camera frame and the ground truth bounding box in the normalized coordinate space transformed by our predicted  $s_{t+1}^{(j)}$ ,  $R_{t+1}^{(j)}$ ,  $T_{t+1}^{(j)}$ . For symmetric objects, we enforce the corner loss on the two intersection points of the symmetric axis and the bounding box surface. The total loss  $L_{\text{total}} = \lambda_{\text{seg}}L_{\text{seg}} + \lambda_{\text{coord}}L_{\text{coord}} + \lambda_{\text{rot}}L_{\text{rot}} + \lambda_{\text{scale}}L_{\text{scale}} + \lambda_{\text{translation}}L_{\text{translation}} + \lambda_{\text{corner}}L_{\text{corner}}$ .

## 5. Experiment

### 5.1. Datasets and Evaluation Metrics

**NOCS-REAL275** For rigid object pose tracking, we evaluate our methods on the NOCS dataset [29] that contains six categories: bottle, bowl, camera, can, laptop, and mug, where bottle, bowl, and can are symmetric. The training set contains: 1) the train split of the CAMERA dataset [29], composed of 300K mixed reality data with synthetic object models from ShapeNetCore [5] as foregrounds and real backgrounds captured in IKEA; and 2) seven real videos capturing challenging cluttered scenes with three object instances in total for each object category. The testing set, NOCS-REAL275, has six real videos de-

Method	NOCS[29]	CASS[6]	CPS++[17]	Oracle ICP	6-PACK[28]	6-PACK [28]	Ours	Ours+RGB seg.
Input	RGBD	RGBD	RGB	Depth	RGBD	RGBD	Depth	RGBD
Setting	Single frame			Tracking				
Initialization	N/A	N/A	N/A	GT.	GT.	Pert.	Pert.	Pert.
$5^\circ 5\text{cm} \uparrow$	16.97	29.44	2.24	0.65	28.92	22.13	62.16	<b>63.60</b>
mIoU $\uparrow$	55.15	55.98	30.02	14.69	55.42	53.58	64.10	<b>69.19</b>
$R_{err} \downarrow$	20.18	14.17	25.32	40.28	19.33	19.66	<b>5.94</b>	6.43
$T_{err} \downarrow$	4.85	12.07	21.62	7.71	<b>3.31</b>	3.62	7.92	4.18

Table 1. **Results of category-level rigid object pose tracking on NOCS-REAL275.** The results are averaged over all 6 categories.

picting in total three different (unseen) instances for each object category totaling 3200 frames.

**Articulated Objects from SAPIEN** For articulated object pose tracking, we create a synthetic dataset using SAPIEN [34]. Our dataset contains four categories: glasses, scissors, laptop, and drawers, where drawers have prismatic joints and the others have revolute joints. The training set contains 98K depth images of 164 standalone object instances with random joint states and viewpoints. The testing set contains 190 depth videos of 19 unseen instances with a total length of 19K frames, where instances keep moving and changing their joint states. See appendix E for more information.

**Real-World Laptop Test Trajectories from the BMVC dataset [18]** We also test our model on real articulated object trajectories from [18], where the objects maintain the same joint state and only viewpoint changes. Among the 4 instances in the dataset, we can only evaluate our method on the laptop for which we have category-level training data from SAPIEN. The two laptop depth sequences contain a total of 1765 frames.

**Evaluation Metrics** We report the following metrics for both rigid and articulated object pose tracking: 1)  **$5^\circ 5\text{cm}$  accuracy**, the percentage of pose predictions with rotation error  $< 5^\circ$  and translation error  $< 5\text{cm}$ ; 2) **mIoU**, the average 3D intersection over union of ground-truth and predicted bounding boxes; 3)  **$R_{err}(\circ)$** , average rotation error; and 4)  **$T_{err}(\text{cm})$** , average translation error. For articulated objects, we additionally report the average joint state accuracy: 5)  **$\theta_{err}(\circ)$**  rotation error for revolute joints; and 6)  **$d_{err}(\text{cm})$**  translation error for prismatic joints. For real-world laptop trajectories, we follow [18] and use pose tolerance, namely the Averaged Distance (AD) accuracy with 10% of the object part diameter as the threshold.

## 5.2. Category-Level Rigid Object Pose Tracking

**Experiment Setting** To track an object in the cluttered scenes from the NOCS-REAL275 dataset, we propose to first crop from the scene a ball of depth points enclosing the object of interest. We set the center and the radius of the ball according to the previous frame’s 9DoF pose estimation. To generate training data, we jitter the ground-truth pose with Gaussian noises ( $\sigma_{scale} = 0.02$ ,  $\sigma_{rot} = 5^\circ$ , and  $\sigma_{trans} = 3\text{cm}$ ) to mimic interframe pose changes and crop balls accordingly. At test time, we randomly sample an initial pose around the ground-truth pose for the first frame

from the same distribution.

**Results** Table 1 summarizes the quantitative results for rigid object pose tracking. We report the performance of our method using only depth and using RGBD where object segmentation masks can be obtained by running off-the-shelf detectors on RGB images as in CASS[6]. We compare our method with: 6-PACK [28], a tracking based method that is initialized with the same pose error distributions or ground-truth poses (6-PACK originally only initializes with translation errors); Oracle ICP, which leverages the ground truth object models; and several single-frame based method, including NOCS [29], CASS [6] and CPS++ [17].

Our method achieves the best performance and significantly outperforms the previous state-of-the-art method, 6-PACK, under both initialization settings. We are especially competitive under the rotation error and  $5^\circ 5\text{cm}$  metrics, showing less than a third of the rotation error and a 105% higher  $5^\circ 5\text{cm}$  percentage compared to 6-PACK. Using only depth, our method generates relatively weaker performance regarding translation error, however, this is only due to the failure to segment out cameras on the real test depth images, given the huge sim2real domain gap between our mostly synthetic training data and noisy real data. See section 5.8 and appendix H.1 for detailed analysis. Excluding this camera category, our method will be the best under all metrics (see appendix H.1). It is worth noting that while our method tracks the full 9DoF pose and predicts the bounding boxes, 6-PACK only tracks the 6DoF rigid transformation and has to use a ground-truth box scale to compute 3D IoU, which unfairly advantages 6-PACK during the comparison. Fig. 1 further shows the qualitative comparison between our method and 6-PACK. Our method loses track less often and gives better pose estimations.

## 5.3. Category-Level Articulated Pose Tracking

In Table 2 and Fig. 1, we show our articulated pose tracking results on held-out test instances from the SAPIEN dataset. We compare our method to 1) ANCSH\* (oracle ANCSH), where we provide ground-truth object segmentation masks to the state-of-the-art single frame articulated object pose estimation work, ANCSH [14]. The original ANCSH fails drastically on part segmentation on our dataset due to the part ambiguity of textureless object point clouds rendered from arbitrary viewpoints; and 2) oracle

Method	5° 5cm↑	mIoU↑	$R_{err}$ ↓	$T_{err}$ ↓	$\theta_{err}$ ↓	$d_{err}$ ↓
ANCSH* [14]	92.55	68.69	2.18	0.48	1.62	0.64
Oracle ICP	62.87	56.61	8.95	3.04	7.21	1.05
Ours	<b>98.35</b>	<b>74.00</b>	<b>1.03</b>	<b>0.29</b>	<b>1.38</b>	<b>0.34</b>
C-sRT regression	21.69	34.21	20.48	11.46	6.08	7.57
C-CoordinateNet	95.06	71.99	2.09	0.40	1.52	0.75
C-Crd. + DSAC++ [3]	95.68	68.21	1.80	0.47	1.61	0.56
Ours w/o $L_c, L_s, L_t$	97.63	72.09	1.24	0.35	1.43	0.36
Ours + Rot. Proj.	98.74	74.17	0.97	0.29	1.37	0.34

Table 2. **Experiment results and ablation studies of articulated object pose tracking on the held-out instances from SAPIEN.**  $\theta_{err}$  is averaged over all revolute joints, while  $d_{err}$  is averaged over all prismatic joints. Other results are averaged over parts and categories. See appendix H.2 for per-part, per-category results. Ours + Rot. Proj leverages kinematic constraints, see Section 5.8.

Method	Michel et al.	ANCSH	ANCSH*	Ours
Setting	Known instance		Category-level	
1 all / parts	64.8 / 65.5 66.9	94.1 / 97.5 94.7	74.7 / 89.1 78.5	<b>95.5 / 99.8 95.7</b>
2 all / parts	65.7 / 66.3 66.6	98.4 / 98.9 <b>99.0</b>	97.0 / 98.0 97.6	<b>98.9 / 100.0 98.9</b>

Table 3. **Results on two real sequences of an unseen laptop** are measured in pose tolerance (the higher, the better, see [18]). The left two columns reported by [18, 14] are directly trained on the instance, whereas ANCSH\*(with GT part mask) and ours are only trained on SAPIEN and have never seen the instance.

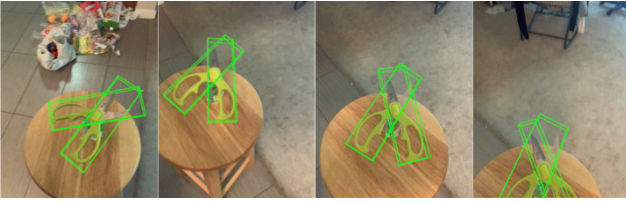


Figure 3. **Qualitative evaluation on the real scissors trajectory.**

ICP, where we assume available ground-truth part labels and object part models, then track each part using ICP.

Note that our articulated SAPIEN dataset is depth-only while RGB-D input is essential for 6-PACK, we thus did not run per-part 6-PACK tracking as a baseline.

We outperform the baselines under all metrics. Although ANCSH\* uses ground-truth labels and regulates its predictions with joint constraints, our per-part scheme still beats it with exceptionally precise rotation estimations.

#### 5.4. Category-Level Articulated Pose Tracking on Real-World Data

We further test our model, trained on the synthetic SAPIEN dataset only, on real-world data. Since the training data does not contain backgrounds, we use pre-segmented object point clouds in the following experiments.

**Real Laptop Trajectories** In Table 3 and Fig.1, we compare our method to Michel et al. [18], ANCSH [14], and ANCSH\* on two real laptop trajectories from [18]. We follow [14] and use their rendered object masks for segmentation. In spite of a Sim2Real gap and a category-level generalization gap, our model outperforms all other methods.

**Real Drawers Under Robot-Object Interaction** To test our method in robotic manipulation scenarios, we capture a

Method	CrdNet	C-Crd.	C-Crd.+ DSAC++	C-sRT	Ours w/o $L_c, L_s, L_t$	Ours
5° 5cm ↑	14.93	46.74	54.77	25.99	60.48	<b>62.16</b>
mIoU↑	49.48	59.99	53.89	32.86	58.80	<b>64.10</b>
$R_{err}$ ↓	53.63	35.08	8.88	34.74	6.41	<b>5.94</b>
$T_{err}$ ↓	9.48	12.97	9.95	21.84	12.64	<b>7.92</b>

Table 4. **Ablation study of rigid object pose tracking on NOCS-REAL275.** The results are averaged over all 6 categories. Here C represents canonicalized.

real drawers trajectory using Kinect2, where a Kinova Jaco2 arm pulls out the middle drawer. As shown in Fig.1(d). Our model successfully tracks both the moving drawer and the other parts. See appendix F for more details.

**Real Scissors Trajectories** Fig.3 shows quantitative results on a real scissors trajectory we captured using Kinect2.

#### 5.5. Ablation Study

To demonstrate the effectiveness of our multi-component design, we construct several variants of our network: 1) CoordinateNet, where we directly regress the NOCS/NPCS coordinates from  $X$  without pose canonicalization; 2) canonicalized CoordinateNet, where we perform pose canonicalization but don't have RotationNet; 3) canonicalized CoordinateNet with DSAC++, where we follow [3] and train our CoordinateNet with a differentiable pose estimation module; 4) canonicalized sRT regression, where we extend our RotationNet to further regress scale and translation based on canonicalized point clouds without using CoordinateNet; and 5) Ours w/o  $L_c, L_s, L_t$  losses, where we discard the pose losses  $L_{scale}, L_{trans}, L_{corner}$  during training. For 1), 2), and 3) we take the coordinate predictions from CoordinateNet and use RANSAC-based pose fitting.

We test the variants on NOCS-REAL275 for rigid object tracking and SAPIEN synthetic dataset for articulated object tracking. The results are summarized in Table 4 and Table 2, where our method outperforms all ablated versions by successfully combining the benefits of pose canonicalization, coordinate prediction, and pose regression. Note that we did not test CoordinateNet without canonicalization on articulated objects due to the part ambiguity of uncolored, arbitrarily posed synthetic objects.

Canonicalized CoordinateNet significantly outperforms CoordinateNet, demonstrating the benefit brought by pose canonicalization. With additional RotationNet, our method further improves canonicalized CoordinateNet and beats the differentiable pipeline CoordinateNet + DSAC++ which also includes explicit pose losses, proving direct regression of small  $\hat{R}_{t+1}^{(j)}$  to be a better choice in the tracking scenario. In contrast, due to ambiguities and insufficient visual cues about scale and translation in the input, the pure regression pipeline, canonicalized sRT regression, produces the worst results. Finally, explicit scale, translation, and corner losses effectively improve our performance compared to ours w/o  $L_c, L_s, L_t$  losses.

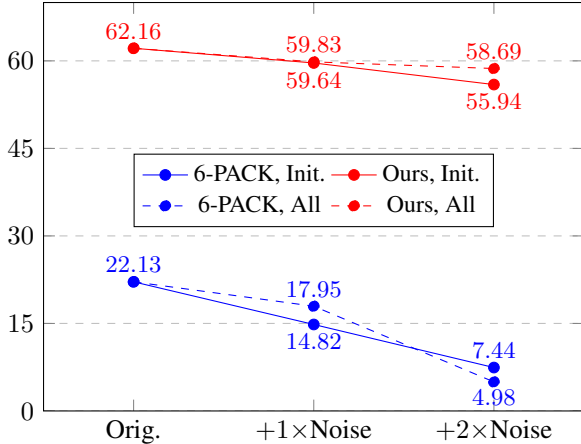


Figure 4.  $5^\circ 5\text{cm}$  (%) w.r.t. **Additional Noise**.  $+m\times\text{Noise}$  means adding  $m$  times train-time errors to 1) the initial pose (denoted **Init.**), which already contains  $1\times$  train-time error; or 2) every frame during training (denoted **All**).

Method	NOCS	6-PACK	ANCSH	Ours
FPS	4.05	3.53	0.80	<b>12.66</b>

Table 5. **Tracking speed in FPS**

## 5.6. Tracking Robustness

Our pose prediction is conditioned on the pose from the previous frame, either the initial pose or an estimated pose. It is therefore worth testing the tracking robustness of our method against noisy pose inputs. As described in Sec. 5.2 and Sec. 5.3, the initial pose errors are randomly sampled from Gaussian distributions. We directly test our model with 1 or 2 times of original pose errors added to (1) the initial pose, and (2) every previous frame’s prediction, to examine the robustness to pose initialization and estimation errors, respectively. We plot the degradation of  $5^\circ 5\text{cm}$  accuracy under each setting and compare to 6-PACK in Fig. 4. Our method is significantly more robust to pose noises than 6-PACK. We are also very robust on articulated objects, see appendix H.3.

## 5.7. Tracking speed

Aside from the state-of-the-art performance, our method also has the highest tracking speed among all others, as summarized in Table 5. All methods are tested on the same device. NOCS and ANCSH are slow due to RANSAC and optimizations, which we don’t need. 6-PACK’s actual speed is slower than what they claim in their paper ( $>10$  FPS) since the network in their officially released code [30] is forwarded 27 times at a grid of potential object centers at each frame to achieve their reported performance.

## 5.8. Discussions

**Tracking Scale** Although the actual scale of the object is constant during tracking, we still track the scale in our framework to deal with inaccurate initial scale. Compared

to fixing the scale as the noisy initial scale throughout the tracking, our scale tracking scheme decreases the average scale error from 1.09% to 0.30% and increases mIoU from 71.70% to 74.00% on articulated objects; and improve mIoU from 73.43% to 76.42% on rigid objects (excluding camera where both schemes fail).

**Leveraging Kinematic Chain Constraints** For articulated objects, our method focuses on per-part tracking without explicitly leveraging joint constraints at test time. Prior works leverage these constraints in instance-level tracking [18, 10] and category-level tracking [14]. [18] and [10] assume perfect knowledge of joint parameters and treat them as a hard constraint. In the category-level setting, however, joint parameters are unknown and difficult to predict due to occlusions, especially for pivot point predictions. Empirically, we have tried predicting them and achieved accuracy similar to state-of-the-art [14], e.g. 1.1cm error for laptop pivot points. However, enforcing these inaccurate constraints harmed our performance, leading to an 80% increase in translation error. ANCSH [14] offers an alternative where only estimated joint axis orientations are used as soft constraints for rotation predictions at the cost of lower speed. Note that without leveraging the constraints, our method already significantly outperforms ANCSH [14]. Without sacrificing speed, we examine the usage of ground truth joint axis orientations as hard constraints but only gain little improvement as shown in Table 2 (Ours + Rot. Proj.). We leave this direction to future works.

**Limitations and Failure Cases** Most of our failure cases come from large sensor noise in real depth images. In extreme cases, e.g. on real cameras from NOCS-REAL275 which are either too reflective or too dark, our CoordinateNet fails to produce reasonable segmentation and the whole pipeline collapses (see appendix H.1). In milder cases, domain gap resulting from sensor noise also degrades our performance. Specific domain adaptation techniques may be needed to deal with this issue, which are beyond the scope of this paper.

## 6. Conclusion

In this paper, for the first time, we tackle the problem of category-level 9DoF pose tracking for both rigid and articulated objects. To achieve this goal, we propose an end-to-end differentiable pose tracking framework consisting of three modules: pose canonicalization, RotationNet, and CoordinateNet. Our algorithm achieves state-of-the-art performance on both category-level rigid and articulated pose benchmarks and runs comparably fast for evaluation.

**Acknowledgement:** This research is supported by a grant from the SAIL-Toyota Center for AI Research, a grant from the Samsung GRO program, NSF grant IIS-1763268, a Vannevar Bush Faculty fellowship, the support of the Stanford UGVR program, and gifts from Kwai and Qualcomm.



## References

- [1] Mykhaylo Andriluka, Umar Iqbal, Eldar Insafutdinov, Leonid Pishchulin, Anton Milan, Juergen Gall, and Bernt Schiele. PoseTrack: A benchmark for human pose estimation and tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5167–5176, 2018. [3](#)
- [2] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. Dsac-differentiable RANSAC for camera localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6684–6692, 2017. [2](#), [3](#)
- [3] Eric Brachmann and Carsten Rother. Learning less is more-6D camera localization via 3D surface regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4654–4662, 2018. [3](#), [7](#)
- [4] Tommaso Cavallari, Stuart Golodetz, Nicholas A Lord, Julien Valentin, Luigi Di Stefano, and Philip HS Torr. On-the-fly adaptation of regression forests for online camera relocalisation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4457–4466, 2017. [2](#)
- [5] 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. [5](#)
- [6] 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 Conference on Computer Vision and Pattern Recognition*, pages 11973–11982, 2020. [2](#), [6](#)
- [7] Xu Chen, Zijian Dong, Jie Song, Andreas Geiger, and Otmar Hilliges. Category level object pose estimation via neural analysis-by-synthesis. In *European Conference on Computer Vision*, pages 139–156. Springer, 2020. [2](#)
- [8] Changhyun Choi and Henrik I Christensen. RGB-D object tracking: A particle filter approach on GPU. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1084–1091. IEEE, 2013. [3](#)
- [9] Xinke Deng, Arsalan Mousavian, Yu Xiang, Fei Xia, Timothy Bretl, and Dieter Fox. PoseRBPF: A rao-blackwellized particle filter for 6D object pose tracking. *Robotics: Science and Systems*, 2019. [3](#)
- [10] Karthik Desingh, Shiyang Lu, Anthony Opipari, and Odest Chadwicke Jenkins. Factored pose estimation of articulated objects using efficient nonparametric belief propagation. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 7221–7227. IEEE, 2019. [8](#)
- [11] Varun Ganapathi, Christian Plagemann, Daphne Koller, and Sebastian Thrun. Real-time human pose tracking from range data. In *Proceedings of the European Conference on Computer Vision*, pages 738–751, 2012. [3](#)
- [12] Shangchen Han, Beibei Liu, Randi Cabezas, Christopher D Twigg, Peizhao Zhang, Jeff Petkau, Tsz-Ho Yu, Chun-Jung Tai, Muzaffer Akbay, Zheng Wang, et al. MEgATrack: monochrome egocentric articulated hand-tracking for virtual reality. *ACM Transactions on Graphics*, 39(4):87–1, 2020. [3](#)
- [13] Alexander Krull, Frank Michel, Eric Brachmann, Stefan Gumhold, Stephan Ihrke, and Carsten Rother. 6-DOF model based tracking via object coordinate regression. In *Asian Conference on Computer Vision*, pages 384–399, 2014. [3](#)
- [14] Xiaolong Li, He Wang, Li Yi, Leonidas J Guibas, A Lynn Abbott, and Shuran Song. Category-level articulated object pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3706–3715, 2020. [2](#), [3](#), [5](#), [6](#), [7](#), [8](#), [11](#)
- [15] 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*, pages 683–698, 2018. [3](#), [4](#)
- [16] Fabian Manhardt, Wadim Kehl, and Adrien Gaidon. Roi-10d: Monocular lifting of 2d detection to 6d pose and metric shape. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2069–2078, 2019. [5](#)
- [17] Fabian Manhardt, Gu Wang, Benjamin Busam, Manuel Nickel, Sven Meier, Luca Minciullo, Xiangyang Ji, and Nassir Navab. Cps++: Improving class-level 6d pose and shape estimation from monocular images with self-supervised learning, 2020. [6](#)
- [18] Frank Michel, Alexander Krull, Eric Brachmann, Michael Ying Yang, Stefan Gumhold, and Carsten Rother. Pose estimation of kinematic chain instances via object coordinate regression. In *BMVC*, pages 181–1, 2015. [1](#), [2](#), [6](#), [7](#), [8](#)
- [19] Maher Moakher. Means and averaging in the group of rotations. *SIAM Journal on Matrix Analysis and Applications*, 24(1):1–16, 2002. [5](#), [11](#)
- [20] Franziska Mueller, Florian Bernard, Oleksandr Sotnychenko, Dushyant Mehta, Srinath Sridhar, Dan Casas, and Christian Theobalt. GANerated hands for real-time 3D hand tracking from monocular RGB. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 49–59, 2018. [3](#)
- [21] Iason Oikonomidis, Nikolaos Kyriazis, and Antonis A Argyros. Efficient model-based 3D tracking of hand articulations using kinect. In *Bmvc*, volume 1, page 3, 2011. [3](#)
- [22] Karl Pauwels and Danica Kragic. Simtrack: A simulation-based framework for scalable real-time object pose detection and tracking. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1300–1307. IEEE, 2015. [3](#)
- [23] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, pages 5099–5108, 2017. [2](#), [5](#), [12](#)
- [24] 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](#)

- [25] Tanner Schmidt, Richard A Newcombe, and Dieter Fox. DART: Dense articulated real-time tracking. In *Robotics: Science and Systems*, 2014. 3
- [26] Meng Tian, Marcelo H Ang, and Gim Hee Lee. Shape prior deformation for categorical 6d object pose and size estimation. In *European Conference on Computer Vision*, pages 530–546. Springer, 2020. 2
- [27] Shinji Umeyama. Least-squares estimation of transformation parameters between two point patterns. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(4):376–380, 1991. 2, 3, 4, 5, 11, 12
- [28] C. Wang, R. Martín-Martín, D. Xu, J. Lv, C. Lu, L. Fei-Fei, S. Savarese, and Y. 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, 2020. 2, 3, 6
- [29] 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, 11, 13
- [30] Jeremy Wang. 6-pack. <https://github.com/j96w/6-PACK>, 2020. 8
- [31] Robert Y Wang and Jovan Popović. Real-time hand-tracking with a color glove. *ACM Transactions on Graphics*, 28(3):1–8, 2009. 3
- [32] Bowen Wen, Chaitanya Mitash, Baozhang Ren, and Kostas E Bekris. se(3)-TrackNet: Data-driven 6D pose tracking by calibrating image residuals in synthetic domains. *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2020. 2, 3, 4
- [33] Manuel Wüthrich, Peter Pastor, Mrinal Kalakrishnan, Jeanette Bohg, and Stefan Schaal. Probabilistic object tracking using a range camera. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3195–3202, 2013. 3
- [34] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and Hao Su. SAPIEN: A simulated part-based interactive environment. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1, 2, 6, 12
- [35] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. In *Robotics: Science and Systems (RSS)*, 2018. 2
- [36] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *Proceedings of the European Conference on Computer Vision*, pages 466–481, 2018. 3
- [37] Li Yi, Haibin Huang, Difan Liu, Evangelos Kalogerakis, Hao Su, and Leonidas Guibas. Deep part induction from articulated object pairs. In *SIGGRAPH Asia 2018 Technical Papers*, page 209, 2018. 5
- [38] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3d: A modern library for 3d data processing. *arXiv preprint arXiv:1801.09847*, 2018. 3
- [39] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5745–5753, 2019. 5, 11