

Category-Level Articulated Object Pose Estimation

Xiaolong Li^{1*} He Wang^{2*} Li Yi³ Leonidas Guibas² A. Lynn Abbott¹ Shuran Song⁴

¹Virginia Tech ²Stanford University ³Google Research ⁴Columbia University

[articulated-pose.github.io](https://github.com/articulated-pose)

Abstract

This paper addresses the task of category-level pose estimation for articulated objects from a single depth image. We present a novel category-level approach that correctly accommodates object instances previously unseen during training. We introduce Articulation-aware Normalized Coordinate Space Hierarchy (ANCSH) – a canonical representation for different articulated objects in a given category. As the key to achieve intra-category generalization, the representation constructs a canonical object space as well as a set of canonical part spaces. The canonical object space normalizes the object orientation, scales and articulations (e.g. joint parameters and states) while each canonical part space further normalizes its part pose and scale. We develop a deep network based on PointNet++ that predicts ANCSH from a single depth point cloud, including part segmentation, normalized coordinates, and joint parameters in the canonical object space. By leveraging the canonicalized joints, we demonstrate: 1) improved performance in part pose and scale estimations using the induced kinematic constraints from joints; 2) high accuracy for joint parameter estimation in camera space.

1. Introduction

Our environment is populated with articulated objects, ranging from furniture such as cabinets and ovens to small tabletop objects such as laptops and eyeglasses. Effectively interacting with these objects requires a detailed understanding of their articulation states and part-level poses. Such understanding is beyond the scope of typical 6D pose estimation algorithms, which have been designed for rigid objects [31, 25, 24, 28]. Algorithms that do consider object articulations [13, 14, 12, 16] often require the exact object CAD model and the associated joint parameters at test time, preventing them from generalizing to new object instances.

In this paper, we adopt a learning-based approach to perform category-level pose estimation for articulated objects. Specifically, we consider the task of estimating per-part 6D poses and 3D scales, joint parameters (i.e. type, position, axis orientation), and joint states (i.e. joint angle)

* indicates equal contributions.

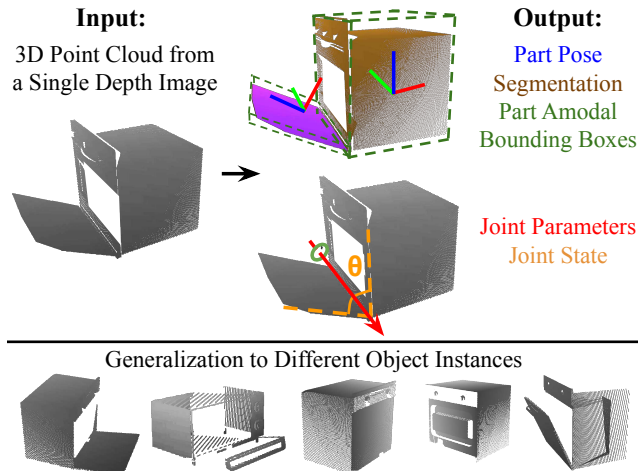


Figure 1. **Category-level articulated object pose estimation.** Given a depth point cloud of a novel articulated object from a known category, our algorithm estimates: part attributes, including part segmentation, poses, scales and amodal bounding boxes; joint attributes, including joint parameters and joint states.

of a novel articulated object instance in a known category from a single depth image. Here object instances from one category will share a known kinematic chain composing of a fixed number of rigid parts connected by certain types of joints. We are particularly interested in the two most common joint types, revolute joints that cause 1D rotational motion (e.g., door hinges), and prismatic joints that allow 1D translational movement (e.g., drawers in a cabinet). An overview is shown in Figure 1. To achieve this goal, several major challenges need to be addressed:

First, to handle novel articulated objects without knowing exact 3D CAD models, we need to find a shared representation for different instances within a given category. The representation needs to accommodate the large variations in part geometry, joint parameters, joint states, and self-occlusion patterns. More importantly, for learning on such diverse data, the representation needs to facilitate intra-category generalization.

Second, in contrast to a rigid object, an articulated object is composed of multiple rigid parts leading to a higher degree of freedom in its pose. Moreover, the parts are connected and constrained by certain joints and hence their

poses are not independent. It is challenging to accurately estimate poses in such a high-dimensional space while complying with physical constraints.

Third, various types of joints provide different physical constraints and priors for part articulations. Designing a framework that can accurately predict the parameters and effectively leverage the constraints for both revolute and prismatic joints is yet an open research problem.

To address the representation challenge, we propose a shared category-level representation for different articulated object instances, namely, Articulation-aware Normalized Coordinate Space Hierarchy (ANCSH). Concretely, ANCSH is a two-level hierarchy of canonical space, composed of Normalized Articulated Object Coordinate Space (NAOCS) at the root level and a set of Normalized Part Coordinate Spaces (NPCSs) at the leaf level. In the NAOCS, object scales, orientations, and joint states are normalized. In the NPCS of each rigid part, the part pose and scale are further normalized. We note that NAOCS and NPCS are complimentary to each other: NAOCS provides a canonical reference on the object level while NPCSs provide canonical part references. The two-level reference frames from ANCSH allow us to define per-part pose as well as joint attributes for previously unseen articulated object instances on the category-level.

To address the pose estimation challenge, we segment objects into multiple rigid parts and predict the normalized coordinates in ANCSH. However, separate per-part pose estimation could easily lead to physically impossible solutions since joint constraints are not considered. To conform with the kinematic constraints introduced by joints, we estimate joint parameters in the NAOCS from the observation, mathematically model the constraints based upon the joint type, and then leverage the kinematic priors to regularize the part poses. We formulate articulated pose fitting from the ANCSH to the depth observation as a combined optimization problem, taking both part pose fitting and joint constraints into consideration. In this work we mainly focus on 1D revolute joints and 1D prismatic joints, while the above formulation can be extended to model and support other types of joints.

Our experiments demonstrate that leveraging the joint constraints in the combined optimization leads to improved performance in part pose and scale prediction. Noting that leveraging joint constraints for regularizing part poses requires high-accuracy joint parameter predictions, which itself is very challenging. Instead of directly predicting joint parameters in the camera space, we consider and leverage predictions in NAOCS, where joints are posed in a canonical orientation, e.g. the revolute joints always point upward for eyeglasses. By transforming joint parameter predictions from NAOCS back to camera space, we further demonstrate supreme accuracy on camera-space

joint parameter predictions. In summary, the primary contribution of our paper is a unified framework for category-level articulated pose estimation. In support of this framework, we design:

- A novel category-level representation for articulated objects – Articulation-aware Normalized Coordinate Space Hierarchy (ANCSH).
- A PointNet++ based neural network that is capable of predicting ANCSH for previously unseen articulated object instances from a single depth input.
- A combined optimization scheme that leverages ANCSH predictions along with induced joint constraints for part pose and scale estimation.
- A two-step approach for high-accuracy joint parameter estimation that first predicts joint parameters in the NAOCS and then transforms them into camera space using part poses.

2. Related Work

This section summarizes related work on pose estimation for rigid and articulated objects.

Rigid object pose estimation. Classically, the goal of pose estimation is to infer an object’s 6D pose (3D rotation and 3D location) relative to a given reference frame. Most previous work has focused on estimating instance-level pose by assuming that exact 3D CAD models are available. For example, traditional algorithms such as iterative closest point (ICP) [4] perform template matching by aligning the CAD model with an observed 3D point cloud. Another family of approaches aim to regress the object coordinates onto its CAD model for each observed object pixel, and then use voting to solve for object pose [6, 7]. These approaches are limited by the need to have exact CAD models for particular object instances.

Category-level pose estimation aims to infer an object’s pose and scale relative to a category-specific canonical representation. Recently, Wang *et al.* [28] extended the object coordinate based approach to perform category-level pose estimation. The key idea behind the intra-category generalization is to regress the coordinates within a Normalized Object Coordinate Space (NOCS), where the sizes are normalized and the orientations are aligned for objects in a given category. Whereas the work by [28] focuses on pose and size estimation for rigid objects, the work presented here extends the NOCS concept to accommodate articulated objects at both part and object level. In addition to pose, our work also infers joint information and addresses particular problems related to occlusion.

Articulated object pose estimation. Most algorithms that attempt pose estimation for articulated objects assume that instance-level information is available. The approaches

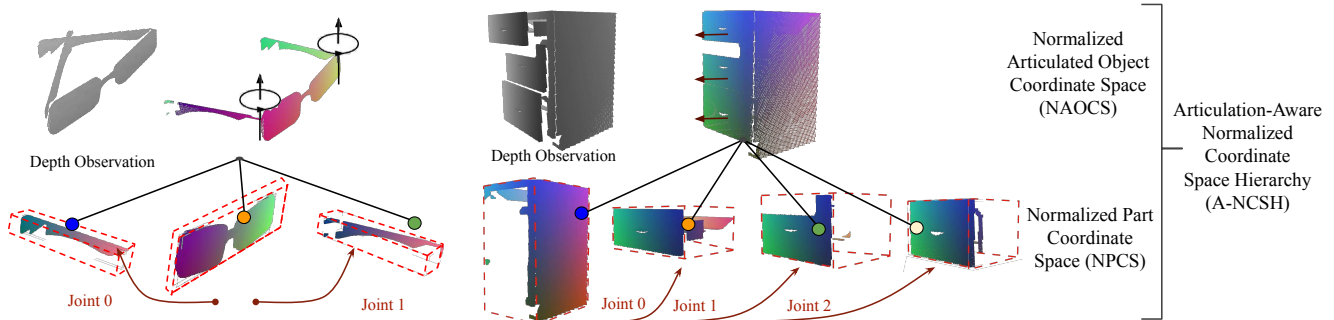


Figure 2. **Articulation-aware Normalized Coordinate Space Hierarchy (ANCSH)** is a category-level object representation composed of a Normalized Articulated Object Coordinate Space (NAOCS) on top of a set of Normalized Part Coordinate Spaces (NPCSs) per part. Here we show two examples of ANCSH representation (points are colored according to its corresponding coordinates in the NAOCS/NPCS). Note that NAOCS sets the object articulations to pre-defined states, all the joints in the NAOCS are hence canonicalized, e.g. the axes of the revolute joints in the eyeglasses example all point upwards and the joint angles are right angles. For each individual part, NPCS maintains the part orientation as in the NAOCS but zero-centers its position and normalizes its scales.

often use CAD models for particular instances along with known kinematic parameters to constrain the search space and to recover the pose separately for different parts [18, 9]. Michel *et al.* [18] used a random forest to vote for pose parameters on canonical body parts for each point in a depth image, followed by a variant of the Kabsch algorithm to estimate joint parameters using RANSAC-based energy minimization. Desingh *et al.* [9] adopted a generative approach using a Markov Random Field formulation, factoring the state as individual parts constrained by their articulation parameters. However, these approaches only consider known object instances and cannot handle different part and kinematic variations. A recent work [1] also tries to handle novel objects within the same category by training a mixed density model [5] on depth images, their method could infer kinematic model using probabilities predictions of a mixtures of Gaussians. However they don't explicitly estimate pose on part-level, the simplified geometry predictions like length, width are for the whole object with scale variation only.

Another line of work relies on active manipulation of an object to infer its articulation pattern [13, 14, 12, 16, 32]. For example, Katz *et al.* [14], uses a robot manipulator to interact with articulated objects as RGB-D videos are recorded. Then the 3D points are clustered into rigid parts according to their motion. Although these approaches could perform pose estimation for unknown objects, they require the input to be a sequence of images that observe an object's different articulation states, whereas our approach is able to perform the task using a single depth observation.

Human body and hand pose estimation. Two specific articulated classes have gained considerable attention recently: the human body and the human hand. For human pose estimation, approaches have been developed using end-to-end networks to predict 3D joint locations directly [17, 23, 19], using dense correspondence maps between 2D images and 3D surface models [3], or

estimating full 3D shape through 2D supervision [15, 20]. Similarly, techniques for hand pose estimation (*e.g.*, [27, 11]) leverages dense coordinate regression, which is then used for voting 3D joint locations. Approaches for both body and hand pose estimation are often specifically customized for those object types, relying on a fixed skeletal model with class-dependent variability (*e.g.*, expected joint lengths) and strong shape priors (*e.g.*, using parametric body shape model for low-dimensional parameterization). Also, such hand/body approaches accommodate only revolute joints. In contrast, our algorithm is designed to handle generic articulated objects with varying kinematic chain, allowing both revolute joints and prismatic joints.

3. Problem Statement

The input to the system is a 3D point cloud $P = \{\mathbf{p}_i \in \mathbb{R}^3 \mid i = 1, \dots, N\}$ backprojected from a single depth image representing an unknown object instance from a known category, where N denotes the number of points. We know that all objects from this category share the same kinematic chain composed of M rigid parts $\{S^{(j)} \mid j = 1, \dots, M\}$ and K joints with known types $\{J_k \mid k = 1, \dots, K\}$. The goal is to segment the point cloud into rigid parts $\{S^{(j)}\}$, recover the 3D rotations $\{R^{(j)}\}$, 3D translations $\{\mathbf{t}^{(j)}\}$, and sizes $\{s^{(j)}\}$ for the parts in $\{S^{(j)}\}$, and predict the joint parameter $\{\phi_k\}$ and state $\{\theta_k\}$ for the joints in $\{J_k\}$. In this work, we consider 1D revolute joints and 1D prismatic joints. We parameterize the two types of joints as following. For a revolute joint, its joint parameters include the direction of the rotation axis $\mathbf{u}_k^{(r)}$ as well as a pivot point \mathbf{q}_k on the rotation axis; its joint state is defined as the relative rotation angle along $\mathbf{u}_k^{(r)}$ between the two connected parts compared with a pre-defined rest state. For a prismatic joint, its joint parameter is the direction of the translation axis $\mathbf{u}_k^{(t)}$, and its joint state is defined as the relative translation distance along $\mathbf{u}_k^{(t)}$ between the two connected parts compared with a pre-defined rest state.

4. Method

ANCSH provides a category-specific reference frame defining per-part poses as well as joint attributes for previously unseen articulated object instances. In Sec. 4.1, we first explain ANCSH in detail. In Sec. 4.2, we then present a deep neural network capable of predicting the ANCSH representation. Sec. 4.3 describes how the ANCSH representation is used to jointly optimize part poses with explicit joint constraints. Last, we describe how we compute joint states and deduce camera-space joint parameters in Sec. 4.4.

4.1. ANCSH Representation

Our ANCSH representation is inspired by and closely related to Normalized Object Coordinate Space (NOCS) [28], which we briefly review here. NOCS is defined as a 3D space contained within a unit cube and was introduced in [28] to estimate the category-level 6D pose and size of rigid objects. For a given category, the objects are consistently aligned by their orientations in the NOCS. Furthermore, these objects are zero-centered and uniformly scaled so that their tight bounding boxes are all centered at the origin of the NOCS with a diagonal length of 1. NOCS provides a reference frame for rigid objects in a given category so that the object pose and size can then be defined using the similarity transformation from the NOCS to the camera space. However, NOCS is limited for representing articulated objects. Instead of the object pose and size, we care more about the poses and the states for each individual parts and joints, which isn't addressed in NOCS.

To define category-level per-part poses and joint attributes, we present ANCSH, a two-level hierarchy of normalized coordinate spaces, as shown in Figure 2. At the root level, NAOCS provides an object-level reference frame with normalized pose, scale, and articulation; at the leaf level, NPCS provides a reference frame for each individual part. We explain both NPCS and NAOCS in detail below.

NAOCS. To construct a category-level object reference frame for the collection of objects, we first bring all the object articulations into a set of pre-defined rest states. Basically, for each joint J_k , we manually define its rest state θ_{k0} and then set the joint into this state. For example, we define the rest states of the two revolute joints in the eyeglasses category to be in right angles; we define the rest states of all drawers to be closed. In addition to normalizing the articulations, NAOCS applies the same normalization used in [28] to the objects, including zero-centering, aligning orientations, and uniformly scaling.

As a canonical object representation, NAOCS has the following advantages: 1) the joints are set to predefined states so that accurately estimating joint parameters in NAOCS, e.g. the direction of rotation/translation axis, becomes an easy task; 2) with the canonical joints, we can

build simple mathematical models to describe the kinematic constraints regarding each individual joint in NAOCS.

NPCS. For each part, NPCS further zero-centers its position and uniformly scales it as is done in [28], while at the same time keeps its orientation unchanged as in NAOCS. In this respect, NPCS is defined similarly to NOCS [28] but for individual parts instead of whole objects. NPCS provides a part reference frame and we can define the part pose and scale as the transformation from NPCS to the camera space. Note that corresponding parts of different object instances are aligned in NPCS, which facilitates intra-category generalization and enables predictions for unseen instances.

Relationship between NPCS, NAOCS and NOCS.

Both NPCS and NAOCS are inspired by the NOCS representation and designed for handling a collection of articulated objects from a given category. Therefore, similar to NOCS, both representations encode canonical information and enable generalization to new object instances. However, each of the two representations has its own advantages in modeling articulated objects and hence provides complementary information. Thus, our ANCSH leverages both NPCS and NAOCS to form a comprehensive representation of both parts and articulations.

On the one hand, NPCSS normalize the position, orientation, and size for each part. Therefore, transformation between NPCSS and camera space can naturally be used to compute per-part 3D amodal bounding boxes, which is not well-presented in NAOCS representation. On the other hand, NAOCS looks at the parts from a holistic view, encoding the canonical relationship of different parts in the object space. NAOCS provides a parent reference frame to those in NPCSS and allows a consistent definition of the joint parameters across different parts. We hence model joints and predict joint parameters in the NAOCS instead of NPCSS. The joint parameters can be used to deduce joint constraints, which can regularize the poses between connected parts. Note that the information defined in NPCS and NAOCS is not mutually exclusive – each NPCS can transform into its counterpart in NAOCS by a uniform scaling and translation. Therefore, instead of independently predicting the full NAOCS representation, our network predicts the scaling and translation parameters for each object part and directly applies it on the corresponding NPCS to obtain the NAOCS estimation.

4.2. ANCSH Network

We devise a deep neural network capable of predicting the ANCSH representation for unseen articulated object instances. As shown in Figure 3, the network takes a depth point cloud P as input and its four heads output rigid part segmentation, dense coordinate predictions in each NPCS,

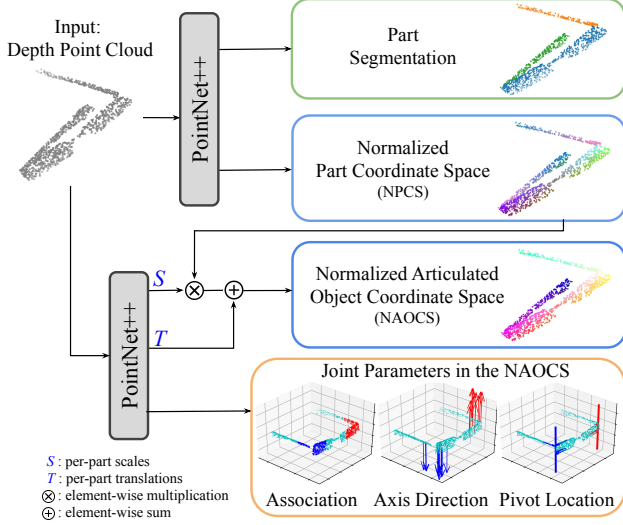


Figure 3. **ANCSH network** leverages two PointNet++ [21] modules to predict the ANCSH representation, including part segmentation, NPCS coordinates, transformations (1D scaling and 3D translation) from each NPCS to the NAOCS, and joint parameters in the NAOCS. This figure illustrates the eyeglasses case with only revolute joints, but the network structure also applies to objects with revolute and prismatic joints.

transformations from each NPCS to NAOCS, and joint parameters in NAOCS, correspondingly. The network is based on two modules adapted from the PointNet++ [21] segmentation architecture.

The part segmentation head predicts a per-point probability distribution among the M rigid parts. The NPCS head predicts M coordinates $\{\mathbf{c}_i^{(j)} \in \mathbb{R}^3 \mid j = 1, \dots, M\}$ for each point \mathbf{p}_i . We use the predicted part label to select the corresponding NPCS. This design helps to inject the geometry prior of each part into the network and hence specializes the networks on part-specific predictions. We design the segmentation network and the NPCS regression network to share the same PointNet++ backbone and only branch at the last fully-connected layers.

The NAOCS head predicts the transformations $\{G^{(j)}\}$ from each NPCS to the NAOCS, and computes the coordinates in NAOCS using the predicted transformations. Since part orientations are the same between NPCS and NAOCS, the network only needs to estimate a 3D translation $G_t^{(j)}$ and a 1D scaling $G_s^{(j)}$ for the NPCS of the part $S^{(j)}$. Similar to NPCS head, the head here predicts for each point \mathbf{p}_i dense transformations with $G_{t,i}^{(j)}$ and $G_{s,i}^{(j)}$ for each NPCS of the parts $S^{(j)}$. We use the predicted segmentation label to select per-point translation $G_{t,i}$ and scaling $G_{s,i}$. Then the NAOCS coordinates can be represented as $\{\mathbf{g}_i \mid \mathbf{g}_i = G_{s,i} \mathbf{c}_i + G_{t,i}\}$. Finally, we compute $G_s^{(j)}$ and $G_t^{(j)}$ by averaging over points $\{\mathbf{p}_i \in S^{(j)}\}$.

The last head infers joint parameters $\{\phi'_k\}$ for each joint

J_k in the NAOCS space (we use “'” here to distinguish the NAOCS parameters from camera-space parameters.) We consider the following two types of joints: 1D revolute joint whose parameters include the rotation axis direction and the pivot point position, namely $\phi'_k = (\mathbf{u}_k^{(r)'}, \mathbf{q}'_k)$; 1D prismatic joint whose parameters is the translation axis direction $\phi'_k = (\mathbf{u}_k^{(t)'})$. We adopt a voting scheme to accurately predict joint parameters, in which we first associate points to each joint via a labeling scheme and then let the points vote for the parameters of its associated joint.

We define a per-point joint association $\{a_i \mid a_i \in \{0, 1, \dots, K\}\}$, where label k means the point \mathbf{p}_i is associated to the joint J_k and label 0 means no association to any joint. We use the following heuristics to provide the ground truth joint association: for a revolute joint J_k , if a point \mathbf{p}_i belongs to its two connecting parts and is within a distance σ from its rotation axis, then we set $a_i = k$; for a prismatic joint, we associate it with all the points on its corresponding moving part. We empirically find $\sigma = 0.2$ leads to a non-overlapping joint association on our data.

In addition to predicting joint association, the joint parameter head performs dense regression on the associated joint parameters. To be more specific, for each point \mathbf{p}_i , the head regresses a 7D vector $\mathbf{v}_i \in \mathbb{R}^7$. The first three dimensions of \mathbf{v}_i is a unit vector, which either represents $\mathbf{u}^{(r)'}$ for a revolute joint or $\mathbf{u}^{(t)'}$ for a prismatic joint. The rest four dimensions are dedicated to the pivot point \mathbf{q}' in case the point is associated to a revolute joint. Since the pivot point of a 1D revolute joint is not uniquely defined (it can move arbitrarily along the rotation axis), we instead predict the projection of \mathbf{p}_i to the rotation axis of its associated revolute joint by regressing a 3D unit vector for the projection direction and a scalar for the projection distance. For training, we only supervise the matched dimensions of \mathbf{v}_i for points \mathbf{p}_i with $a_i \neq 0$. We use the ground truth joint parameters ϕ'_{a_i} associated with joint J_{a_i} as the supervision. During inference, we use the predicted joint association to interpret \mathbf{v}_i . We perform a voting step to get the final joint parameter prediction ϕ'_k , where we simply average the predictions from points associated with each joint J_k . Note that the NAOCS head and the joint parameter head share the second PointNet++ as their backbone since they all predict attributes in the NAOCS.

Loss functions: We use relaxed IoU loss [32] L_{seg} for part segmentation as well as for joint association $L_{\text{association}}$. We use mean-square loss L_{NPCS} for NPCS coordinate regression. We use mean-square loss L_{NAOCS} for NAOCS to supervise per-point translation $\{G_{t,i}^{(j)}\}_{i,j}$ and scaling $\{G_{s,i}^{(j)}\}_{i,j}$. We again use mean-square loss L_{joint} for joint parameter predictions. Our total loss is given by $L = \lambda_1 L_{\text{seg}} + \lambda_2 L_{\text{NPCS}} + \lambda_3 L_{\text{NAOCS}} + \lambda_4 L_{\text{association}} + \lambda_5 L_{\text{joint}}$, where the loss weights are set to $[1, 10, 1, 1, 1]$.

4.3. Pose Optimization with Kinematic Constraints

Given the output of our ANCSH network, including part segmentation, $\{\mathbf{c}_i\}$ for each point \mathbf{p}_i , $\{G_t^{(j)}, G_s^{(j)}\}$ for each part $S^{(j)}$, and $\{\phi'_k\}$ for each joint J_k , we now estimate the 6D poses and sizes $\{R^{(j)}, \mathbf{t}^{(j)}, s^{(j)}\}$ for each part $S^{(j)}$.

Considering a part $S^{(j)}$, for the points $\{\mathbf{p}_i \in S^{(j)}\}$, we have their corresponding NPCS predictions $\{\mathbf{c}_i | \mathbf{p}_i \in S^{(j)}\}$. We could follow [28] to perform pose fitting, where the Umeyama algorithm [26] is adopted within a RANSAC [10] framework to robustly estimate the 6D pose and size of a single rigid object. However, without leveraging joint constraints, naively applying this approach to each individual part in our setting would easily lead to physically impossible part poses. To cope with this issue, we propose the following optimization scheme leveraging kinematic constraints for estimating the part poses. Without the kinematic constraints, the energy function E_{vanilla} regarding all part poses can be written as $E_{\text{vanilla}} = \sum_j e_j$, where

$$e_j = \frac{1}{|S^{(j)}|} \sum_{\mathbf{p}_i \in S^{(j)}} \|\mathbf{p}_i - (s^{(j)} R^{(j)} \mathbf{c}_i + \mathbf{t}^{(j)})\|^2$$

We then introduce the kinematic constraints by adding an energy term e_k for each joint to the energy function. In concrete terms, our modified energy function is $E_{\text{constrained}} = \sum_j e_j + \lambda \sum_k e_k$, where e_k is defined differently for each type of joint. For a revolute joint J_k with parameters $\phi'_k = (\mathbf{u}_k^{(r)'}, \mathbf{q}'_k)$ in the NAOCS, assuming it connects part $S^{(j_1)}$ and part $S^{(j_2)}$, we define e_k as:

$$e_k = \|R^{(j_1)} \mathbf{u}_k^{(r)'} - R^{(j_2)} \mathbf{u}_k^{(r)'}\|^2$$

For a prismatic joint J_k with parameters $\phi'_k = (\mathbf{u}_k^{(t)'})$ in the NAOCS, again assuming it connects part $S^{(j_1)}$ and part $S^{(j_2)}$, we define e_k as:

$$e_k = \mu \|R^{(j_1)} R^{(j_2)T} - I\|^2 + \sum_{j=j_1, j_2} \| [R^{(j)} \mathbf{u}_k^{(t)'}]_{\times} \delta_{j_1, j_2} \|^2$$

where $[\cdot]_{\times}$ converts a vector into a matrix for conducting cross product with other vectors, and δ_{j_1, j_2} is defined as:

$$\delta_{j_1, j_2} = \mathbf{t}^{(j_2)} - \mathbf{t}^{(j_1)} + s^{(j_1)} R^{(j_1)} G_t^{(j_1)} - s^{(j_2)} R^{(j_2)} G_t^{(j_2)}$$

To minimize our energy function $E_{\text{constrained}}$, we can no longer separately solve different part poses using the Umeyama algorithm. Instead, we first minimize E_{vanilla} using the Umeyama algorithm to initialize our estimation of the part poses. Then we fix $\{s^{(j)}\}$ and adopt a non-linear least-squares solver to further optimize $\{R^{(j)}, \mathbf{t}^{(j)}\}$, as is commonly done for bundle adjustment [2]. Similar to [28], we also use RANSAC for outlier removal.

Finally, for each part $S^{(j)}$, we use the fitted $R^{(j)}, \mathbf{t}^{(j)}, s^{(j)}$ and the NPCS $\{\mathbf{c}_i | \mathbf{p}_i \in S^{(j)}\}$ to compute an amodal bounding box, the same as in [28].

4.4. Camera-Space Joint Parameters and Joint States Estimation

Knowing $\{R^{(j)}, \mathbf{t}^{(j)}, s^{(j)}, G_t^{(j)}, G_s^{(j)}\}$ of each part, we can compute the joint states $\{\theta_k\}$ and deduce joint parameters $\{\phi_k\}$ in the camera space from NAOCS joint parameters $\{\phi'_k\}$. For a revolute joint J_k connecting parts $S^{(j_1)}$ and $S^{(j_2)}$, we compute its parameters $\phi_k = (\mathbf{u}_k^{(r)}, \mathbf{q}_k)$ in the camera space as:

$$\mathbf{u}_k^{(r)} = \frac{(R^{(j_1)} + R^{(j_2)}) \mathbf{u}_k^{(r)'}}{\|(R^{(j_1)} + R^{(j_2)}) \mathbf{u}_k^{(r)'}\|}$$

$$\mathbf{q}_k = \frac{1}{2} \sum_{j=j_1, j_2} \frac{R^{(j)} s^{(j)}}{G_s^{(j)}} (\mathbf{q}'_k - G_t^{(j)}) + \mathbf{t}^{(j)}$$

The joint state θ_k can be computed as:

$$\theta_k = \arccos((\text{trace}(R^{(j_2)} (R^{(j_1)})^T) - 1)/2)$$

For a prismatic joint J_k connecting parts $S^{(k_1)}$ and $S^{(k_2)}$, we compute its parameters $\phi_k = (\mathbf{u}_k^{(t)})$ in the camera space similar to computing $\mathbf{u}_k^{(r)}$ for revolute joints and its state θ_k is simply $\|\delta_{k_1, k_2}\|$.

5. Evaluation

5.1. Experimental Setup

Evaluation Metrics. We use the following metrics to evaluate our method.

- **Part-based metrics.** For each part, we evaluate rotation error measured in degrees, translation error, and 3D intersection over union (IoU) [22] of the predicted amodal bounding box.
- **Joint states.** For each revolute joint, we evaluate joint angle error in degrees. For each prismatic joint, we evaluate the error of relative translation amounts.
- **Joint parameters.** For each revolute joint, we evaluate the orientation error of the rotation axis in degrees, and the position error using the minimum line-to-line distance. For each prismatic joint, we compute the orientation error of the translation axis.

Datasets. We have evaluated our algorithm using both synthetic and real-word datasets. To generate the synthetic data, we mainly use object CAD models from [29] along with drawer models from [30]. Following the same rendering pipeline with random camera viewpoints, we use PyBullet[8] to generate on average 3000 testing images of unseen object instances for each object category that do not overlap with our training data. For the real data, we evaluated our algorithm on the dataset provided by Michel *et al.* [18], which contains depth images for 4 different objects captured using the Kinect.

Category	Method	Part-based Metrics			Joint States	Joint Parameters	
		Rotation Error ↓	Translation Error ↓	3D IoU % ↑	Error ↓	Angle error ↓	Distance error ↓
Eye-glasses	NPCS	4.0°, 7.7°, 7.2°	0.044, 0.080, 0.071	86.9, 40.5, 41.4	8.8°, 8.4°	-	-
	NAOCS	4.2°, 12.1°, 13.5°	0.157, 0.252, 0.168	-	13.7°, 15.1°	-	-
	ANCSH	3.7°, 5.1°, 3.7°	0.035, 0.051, 0.057	87.4, 43.6, 44.5	4.3°, 4.5°	2.2°, 2.3°	0.019, 0.014
Oven	NPCS	1.3°, 3.5°	0.032, 0.049	75.8, 88.5	4.0°	-	-
	NAOCS	1.7°, 4.7°	0.036, 0.090	-	5.1°	-	-
	ANCSH	1.1°, 2.2°	0.030, 0.046	75.9, 89.0	2.1°	0.8°	0.024
Washing Machine	NPCS	1.1°, 2.0°	0.043, 0.056	86.9, 88.0	2.3°	-	-
	NAOCS	1.1°, 3.3°	0.072, 0.119	-	3.1°	-	-
	ANCSH	1.0°, 1.4°	0.042, 0.053	87.0, 88.3	1.00°	0.7°	0.008
Laptop	NPCS	11.6°, 4.4°	0.098, 0.044	35.7, 93.6	14.4°	-	-
	NAOCS	12.4°, 4.9°	0.110, 0.049	-	15.2°	-	-
	ANCSH	6.7°, 4.3°	0.062, 0.044	41.1, 93.0	9.7°	0.5°	0.017
Drawer	NPCS	1.9°, 3.5°, 2.4°, 1.8°	0.032, 0.038, 0.024, 0.025	82.8, 71.2, 71.5, 79.3	0.026, 0.031, 0.046	-	-
	NAOCS	1.5°, 2.5°, 2.5°, 2.0°	0.044, 0.045, 0.073, 0.054	-	0.043, 0.066, 0.048	-	-
	ANCSH	1.0°, 1.1°, 1.2°, 1.5°	0.024, 0.021, 0.021, 0.033	84.0, 72.1, 71.7, 78.6	0.011, 0.020, 0.030	0.8°, 0.8°, 0.8°	-

Table 1. **Performance comparison on unseen object instances.** The categories eyeglasses, oven, washing machine, and laptop contain only revolute joints and the drawer category contains three prismatic joints.

Baselines. There are no existing methods for category-level articulated object pose estimation. We therefore use ablated versions of our system for baseline comparison.

- **NPCS.** This algorithm predicts part segmentation and NPCS for each part (without the joint parameters). The prediction allows the algorithm to infer part pose, amodal bounding box for each part, and joint state for revolute joint by treating each part as an independent rigid body. However, it is not able to perform a combined optimization with the kinematic constraints.
- **NAOCS.** This algorithm predicts part segmentation and NAOCS representation for the whole object instance. The prediction allows the algorithm to infer part pose and joint state, but not the amodal bounding boxes for each part since the amodal bounding boxes are not defined in the NAOCS alone. Note the part pose here is defined from the NAOCS to the camera space, different from the one we defined based upon NPCS. We measure the error in the observed object scale so that it is comparable with our method.
- **Direct joint voting.** This algorithm directly votes for joint-associated parameters in camera space, including offset vectors and orientation for each joint from the point cloud using PointNet++ segmentation network.

Our final algorithm predicts the full ANCSH representation that includes NPCS, joint parameters, and per-point global scaling and translation value that can be used together with the NPCS prediction for computing NAOCS.

5.2. Experimental Results

Figure 4 presents some qualitative results. Tables 1 summarizes the quantitative results. Following paragraphs provide our analysis and discussion of the results.

Effect of combined optimization. First, we want to examine how combined optimization would influence the accuracy of articulated object pose estimation, using both predicted joint parameters and predicted part poses. To see this, we compare the algorithm performance between NPCS and ANCSH, where NPCS performs a per-part pose estimation and ANCSH performs a combined optimization using the full kinematic chain to constrain the result. The results in Table 1 show that the combined optimization of joint parameters and part pose consistently improves the predict results for almost all object categories and on almost all evaluation metrics. The improvement is particularly salient for thin object parts such as the two temples of eyeglasses (the parts that extend over the ears), where the per-part based method produces large pose errors due to limited number of visible points and shape ambiguity. This result demonstrates that the joint parameters predicted in the NAOCS can regularize the part poses based on kinematic chain constraints during the combined pose optimization step and improve the pose estimation accuracy.

Joint parameters estimation. Predicting the location and the orientation of joints in camera space directly with all degrees of freedom is challenging. Our approach predicts the joint parameters in NAOCS since it provides a canonical representation where the joint axes usually have a strong orientation prior. We further use a voting-based scheme to reduce the prediction noise. Given joint axis predictions in NAOCS, we leverage the transformation between NAOCS and NPCS to compute corresponding joint parameters in NPCS. Based on the high-quality prediction of part poses, we will transform the joint parameters into the camera coordinate. Comparing to a direct voting baseline using PointNet++, our approach significantly improves the joint axis prediction for unseen instances (Table 2).

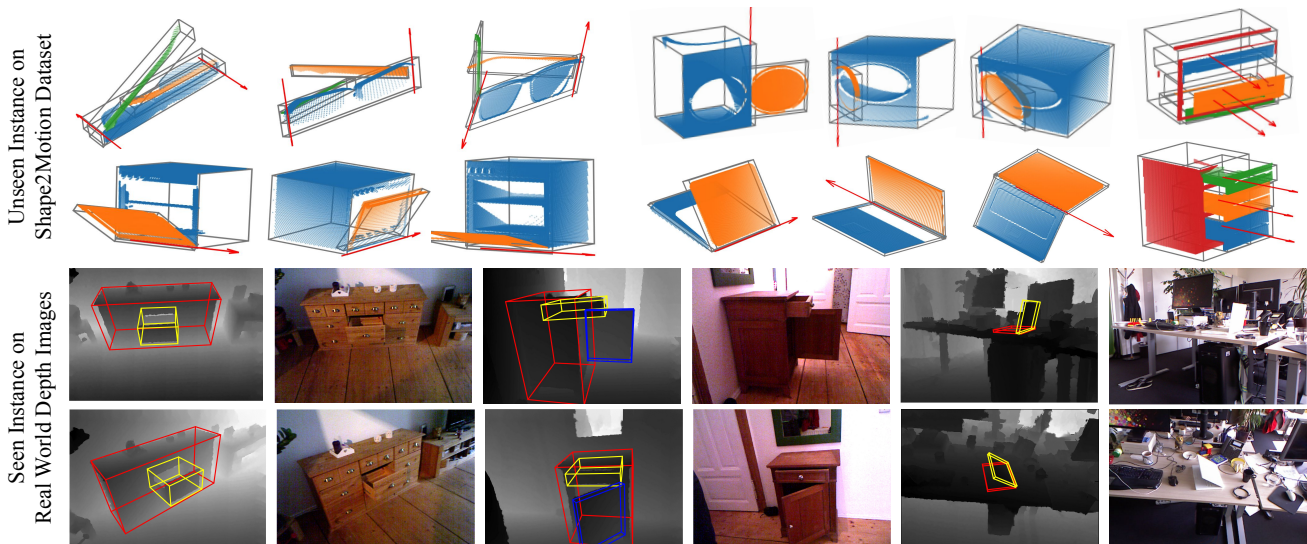


Figure 4. **Qualitative Results.** Top two rows show test results on unseen object instances from the Shape2Motion dataset [29] and SAPIEN dataset [30] (for only drawer category). Bottom two rows show test result on seen instances in the real-world dataset [18]. Here we visualize the predicted amodal bounding box for each part. Color images are for visualization only.

Category	Methods	Angle error	Distance error
Eye-glass	PointNet++	2.9°, 15.7°	0.140, 0.197
	ANCSH	2.2°, 2.3°	0.019, 0.014
Oven	PointNet++	27.0°	0.024
	ANCSH	0.8°	0.024
Washing Machine	PointNet++	8.7°	0.010
	ANCSH	0.7°	0.008
Laptop	PointNet++	29.5°	0.007
	ANCSH	0.5°	0.017
Drawer	PointNet++	4.9°, 5.0°, 5.1°	-
	ANCSH	0.8°, 0.8°, 0.8°	-

Table 2. **A comparison of joint parameters estimation.** Here PointNet++ denotes the direct joint voting baseline.

Generalization to real depth images. We have also tested our algorithm’s ability to generalize to real-world depth images on the dataset provided in [18]. The dataset contains video sequences captured with Kinect for four different object instances. Following the same training protocol, we train the algorithm with synthetically rendered depth images of the provided object instances. Then we test the pose estimation accuracy on the real world depth images. We adopt the same evaluation metric in [18], which uses 10% of the object part diameter as the threshold to compute Averaged Distance (AD) accuracy, and test the performance on each sequence separately. Although our algorithm is not specifically designed for instance-level pose estimation and the network has never been trained using any real-world depth images, our algorithm achieves strong performance on par with or even better

than state-of-the-art. On average our algorithm achieves 96.25%, 92.3%, 96.9%, 79.8% AD accuracy on the whole kinematic chain of object instance laptop, cabinet, cupboard and toy train. For detailed results on each part in all the test sequences, as well as more visualizations, please refer to the supplementary material.

6. Conclusion

This paper has presented an approach for category-level pose estimation of articulated objects from a single depth image. To accommodate unseen object instances with large intra-category variations, we introduce a novel object representation, namely Articulation-aware Normalized Coordinate Space Hierarchy (ANCSH). We further devise a deep neural network capable of predicting ANCSH from a single depth point cloud. We then formulate articulated pose fitting from the ANCSH predictions as a combined optimization problem, taking both part pose errors and joint constraints into consideration. Our experiments demonstrate that the ANCSH representation and the combined optimization scheme significantly improve the accuracy for both part pose prediction and joint parameters estimation.

Acknowledgement: This research is supported by a grant from Toyota-Stanford Center for AI Research, resources provided by Advanced Research Computing in the Division of Information Technology at Virginia Tech. We thank Vision and Learning Lab at Virginia Tech for help on visualization tool. We are also grateful for the financial and hardware support from Google.

References

- [1] Ben Abbatematteo, Stefanie Tellex, and George Konidaris. Learning to generalize kinematic models to novel objects. In *Proceedings of the Third Conference on Robot Learning*, 2019. 3
- [2] Sameer Agarwal, Noah Snavely, Steven M Seitz, and Richard Szeliski. Bundle adjustment in the large. In *European conference on computer vision*, pages 29–42. Springer, 2010. 6
- [3] Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7297–7306, 2018. 3
- [4] Paul J Besl and Neil D McKay. A method for registration of 3-d shapes. In *PAMI*, 1992. 2
- [5] Christopher M Bishop. Mixture density networks. 1994. 3
- [6] Eric Brachmann, Alexander Krull, Frank Michel, Stefan Gumhold, Jamie Shotton, and Carsten Rother. Learning 6d object pose estimation using 3d object coordinates. In *European conference on computer vision*, pages 536–551. Springer, 2014. 2
- [7] Eric Brachmann, Frank Michel, Alexander Krull, Michael Ying Yang, Stefan Gumhold, et al. Uncertainty-driven 6d pose estimation of objects and scenes from a single rgb image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3364–3372, 2016. 2
- [8] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. <http://pybullet.org>, 2016–2018. 6
- [9] Karthik Desingh, Shiyang Lu, Anthony Opipari, and Odest Chadwicke Jenkins. Factored pose estimation of articulated objects using efficient nonparametric belief propagation. *arXiv preprint arXiv:1812.03647*, 2018. 3
- [10] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981. 6
- [11] Lihao Ge, Zhou Ren, and Junsong Yuan. Point-to-point regression pointnet for 3d hand pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 475–491, 2018. 3
- [12] Karol Hausman, Scott Niekum, Sarah Osentoski, and Gaurav S Sukhatme. Active articulation model estimation through interactive perception. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3305–3312. IEEE, 2015. 1, 3
- [13] Dov Katz and Oliver Brock. Manipulating articulated objects with interactive perception. In *2008 IEEE International Conference on Robotics and Automation*, pages 272–277. IEEE, 2008. 1, 3
- [14] Dov Katz, Moslem Kazemi, J Andrew Bagnell, and Anthony Stentz. Interactive segmentation, tracking, and kinematic modeling of unknown 3d articulated objects. In *2013 IEEE International Conference on Robotics and Automation*, pages 5003–5010. IEEE, 2013. 1, 3
- [15] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J Black, and Peter V Gehler. Unite the people: Closing the loop between 3d and 2d human representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6050–6059, 2017. 3
- [16] Roberto Martín-Martín, Sebastian Höfer, and Oliver Brock. An integrated approach to visual perception of articulated objects. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5091–5097. IEEE, 2016. 1, 3
- [17] Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. Vnect: Real-time 3d human pose estimation with a single rgb camera. *ACM Transactions on Graphics (TOG)*, 36(4):44, 2017. 3
- [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. 3, 6, 8
- [19] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Coarse-to-fine volumetric prediction for single-image 3d human pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7025–7034, 2017. 3
- [20] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 459–468, 2018. 3
- [21] 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. 5

- [22] Shuran Song and Jianxiong Xiao. Deep sliding shapes for amodal 3d object detection in rgb-d images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 808–816, 2016. [6](#)
- [23] Xiao Sun, Jiayang Shang, Shuang Liang, and Yichen Wei. Compositional human pose regression. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2602–2611, 2017. [3](#)
- [24] Martin Sundermeyer, Zoltan-Csaba Marton, Maximilian Durner, Manuel Brucker, and Rudolph Triebel. Implicit 3d orientation learning for 6d object detection from rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 699–715, 2018. [1](#)
- [25] 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](#)
- [26] Shinji Umeyama. Least-squares estimation of transformation parameters between two point patterns. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (4):376–380, 1991. [6](#)
- [27] Chengde Wan, Thomas Probst, Luc Van Gool, and Angela Yao. Dense 3d regression for hand pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5147–5156, 2018. [3](#)
- [28] 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](#), [4](#), [6](#)
- [29] Xiaogang Wang, Bin Zhou, Yahao Shi, Xiaowu Chen, Qinqing Zhao, and Kai Xu. Shape2motion: Joint analysis of motion parts and attributes from 3d shapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8876–8884, 2019. [6](#), [8](#)
- [30] 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, 2020. [6](#), [8](#)
- [31] 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. [1](#)
- [32] 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. ACM, 2018. [3](#), [5](#)