

A Hand Motion-guided Articulation and Segmentation Estimation

Richard Sahala Hartanto, Ryoichi Ishikawa, Menandro Roxas, Takeshi Oishi

ABSTRACT

In this paper, we present a method for simultaneous articulation model estimation and segmentation of an articulated object in RGB-D images using human hand motion. Our method uses the hand motion in the processes of the initial articulation model estimation, ICP-based model parameter optimization, and region selection of the target object. The hand motion gives an initial guess of the articulation model: prismatic or revolute joint. The method estimates the joint parameters by aligning the RGB-D images with the constraint of the hand motion. Finally, the target regions are selected from the cluster regions which move symmetrically along with the articulation model. Our experimental results show the robustness of the proposed method for the various objects.

I. INTRODUCTION

Understanding the articulation of objects is crucial for robots to manipulate various functional objects. The robot needs to learn how to use the tools that humans use at home or in the workplace. For that, the robot needs to understand movement patterns and ranges in 3D space, the articulation model, of these objects to manipulate them correctly and to avoid damages of the object, environment, and robot itself.

Most of the previous works on articulation estimation utilize visual input: an RGB or RGB-D image sequence. The image-based method tracks feature points of a moving object in a 2D space [1], [2], [3] or a 3D space [4], [5]. If depth information is given, 3D geometric features are useful [6], [7]. Recently, Deep Neural Network-based approaches such as CNN for functional estimation [8] or a verbal-guided approach [9] have also been proposed. Aside from visually supervised methods, the interactive approach using a robot arm is an alternative research direction [10], but we consider the visual supervision for general-purposes.

In addition to the articulation model estimation, the functional segmentation of the object is also important. One of the popular approaches is a structure-driven method that assumes the target object is composed of rectangular planes [11]. The CNN-based approach is also successful in the segmentation process [8]. Another stream is more closely linked to the articulation, *i.e.*, tracking and clustering movements. The point clouds that belong to a functional segment move along with the same articulation model. Therefore, the segmentation is performed by tracking 2D feature points [1], [12] or

Richard Sahala Hartanto, Ryoichi Ishikawa, Menandro Roxas and Takeshi Oishi are with Institute of Industrial Science, The University of Tokyo, Japan {richard, ishikawa, roxas, oishi}@cvl.iis.u-tokyo.ac.jp

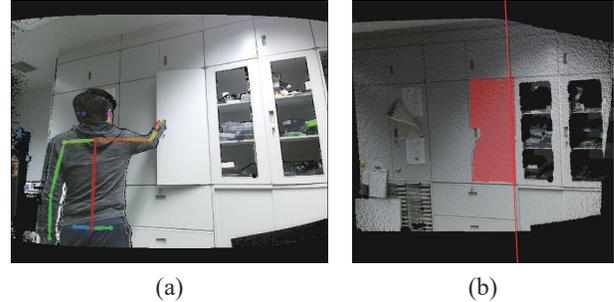


Fig. 1. Articulation and segmentation estimation from RGBD images with using human's hand movement. (a) input RGB-D image sequence with human pose detection by OpenPose [15]. (b) Articulation and segmentation estimation results.

3D structural feature points [13], [14] and by clustering the motions of the points.

Some problems with conventional methods are the dependency on textures and distinctive shapes of the object. Since functional objects are usually textureless and symmetric, these problems are highlighted especially when using noisy depth information from optical sensors. In short, robustly tracking the target regions in 3D space is no easy task.

We consider using human motion while manipulating the articulated object. The human interacts with the target object, and the contact point moves according to the articulation model. Recent CNN development makes it easy and accurate to detect human motion from RGB or RGB-D images [16], [17], [15], even in situations with partial occlusions. We can easily observe the human motion during manipulation and the target regions at the same time.

In this paper, we propose a method of hand motion-guided articulation and segmentation estimation from RGB-D images. The inputs are the RGB-D images of a background scene and a sequence of manipulation (Fig. 1 (a)). The outputs are the articulation model and segmentation of the target object (Fig. 1 (b)). We detect hand motion using a CNN-based method [15] and utilize it for initial estimation of the articulation model, weight and constraint for depth image alignment, and region selection. The proposed method does not primarily depend on the object shape and textures since we estimate the initial model from the human body information, which is characteristic and detectable using modern CNN-based methods. The experimental results show the robustness of our method for various objects. The code

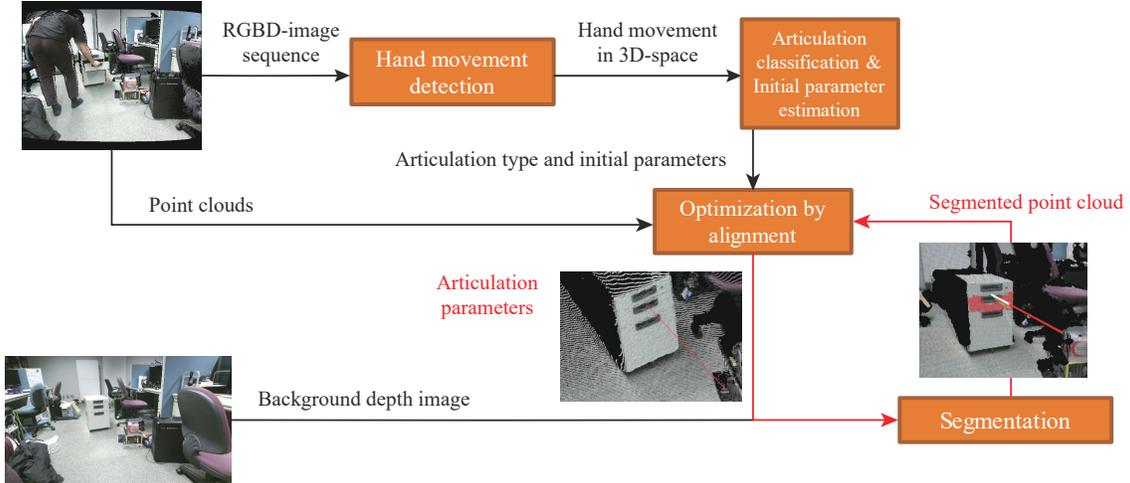


Fig. 2. Overview of our method. The red letter indicates the outputs of our method. The loop indicated by the red arrow is the iterative refinement process.

is released as an open-source¹.

II. OVERVIEW AND NOTATION

A. Overview

We use an RGB-D video that captures a scene of manipulating the articulated object and an RGB-D image of background. Figure 2 shows an overview of the proposed method.

We assume the following conditions:

- Input RGB-D video includes only a manipulation scene, and the first frame is the start of manipulation.
- Only one person is in the RGB-D video.
- It is known which hand is used.
- Articulation type is either prismatic or revolute as described in II-B.

First, we estimate the hand motion in 3D space from RGB images and depth images using CNN-based semantic segmentation and human pose estimation method. Next, the articulation type and the initial parameters are estimated from the hand motion. The alignment is applied to the point clouds derived from the depth images to optimize the articulation parameters according to the estimated articulation type. The segmentation is performed by extracting points which are symmetric and move according to the articulation model. The alignment using the segmentation result refines the parameters. The refinement process is shown as the red arrows in Fig. 2. (In our experiments, we iterate the refinement process twice)

B. Articulation model

We treat an object as either containing a prismatic joint or a revolute joint. These joints are the basis of most objects

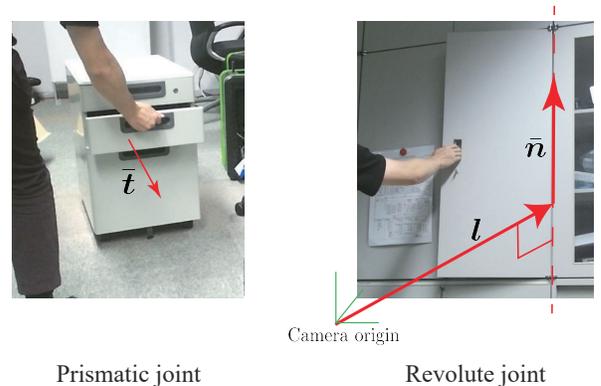


Fig. 3. Example of a prismatic joint and a revolute joint. The prismatic joint object transformed along the direction \bar{t} . The revolute joint object rotates around the axis where going through l in the direction \bar{n} .

designed for humans. We describe how we express these joints as follows.

1) *Prismatic joint*: When a rigid object is linked to a sliding linkage, in which we call the prismatic joint, the movement of the object is restricted to make only a linear sliding movement relative to the linkage (See the left of Fig. 3). The prismatic joint has only one direction to move in 3D space. We express the prismatic joint as 2-DoF normalized 3×1 vector \bar{t} indicating moving direction. The amount of movement of a prismatic object is represented by a . The unit of a can be any but is a meter in the experiment. An object moving along the prismatic joint is transformed by the following 4×4 matrix;

$$\begin{pmatrix} \mathbf{I}_{3 \times 3} & a\bar{t} \\ \mathbf{0}_{1 \times 3} & 0 \end{pmatrix}. \quad (1)$$

¹<https://github.com/cln515/Articulation-Estimation>

where, $\mathbf{I}_{3 \times 3}$ is 3×3 identity matrix.

2) *Revolute joint*: When a rigid object is linked to a rotating linkage, in which we call the revolute joint, the movement of the object is restricted to only rotation around an axis with a certain movement range relative to the linkage (See the right of Fig. 3). The revolute joint can be expressed as a line in 3D space. We express the revolute joint as a total of 4-DoF parameters: 2-DoF unit vector $\bar{\mathbf{n}}$ of the rotation axis and 2-DoF vector \mathbf{l} where the rotation axis passes through and holds $\bar{\mathbf{n}} \cdot \mathbf{l} = 0$. The amount of movement of a revolute object is represented by θ , and the unit is the radian. An object moving along the revolute joint is transformed by the following 4×4 matrix;

$$\begin{pmatrix} \mathbf{R}(\theta, \bar{\mathbf{n}}) & \mathbf{l} - \mathbf{R}(\theta, \bar{\mathbf{n}})\mathbf{l} \\ \mathbf{0}_{1 \times 3} & 0 \end{pmatrix}. \quad (2)$$

where, $\mathbf{R}(\theta, \bar{\mathbf{n}})$ is 3×3 rotation matrix indicating rotation by θ around $\bar{\mathbf{n}}$:

$$\mathbf{R}(\theta, \bar{\mathbf{n}}) = \cos \theta \mathbf{I} + (1 - \cos \theta) \bar{\mathbf{n}} \bar{\mathbf{n}}^\top + \sin \theta [\bar{\mathbf{n}}]_\times, \quad (3)$$

where $[\cdot]_\times$ is 3×3 skew symmetric matrix.

III. METHODOLOGY

A. Initial articulation estimation

We recognize the articulation type and estimate the initial articulation parameters from the hand motion during manipulation in the 3D space. To detect the hand position in RGB images, we use hand keypoint detection [18] implemented in OpenPose [15]. The hand keypoint detection offers a detailed hand pose as a set of 2D points and it needs a bounding box around the hand area. The bounding box is offered from OpenPose when the entire body is pictured or CNN-based hand detection implementation [19] trained using the EgoHands dataset [20]. After obtaining a hand pose in 2D space, we compute the 3D points of the hand points using the depth map and compute the 3D centroid of the hand.

After estimating the hand motion, we fit a circle to the motion path for recognizing the articulation type. We use RANSAC and a non-linear optimization method to fit the circle to the noisy points. Counting the number of inlier points is not practical since a straight line can be considered as a part of a large circle. Therefore, we use the range of the movement angle for type recognition. If the range is smaller than a predetermined threshold value (we set the value to 30 degrees), we recognize the target object to be prismatic and perform the line fitting.

We estimate the initial articulation parameters from the result of the fitting. The line fitting computes 4-DoF parameters as explained in Sec. II-B.1. The circle parameters are 6-DoF parameters consisting of the radius r (1-DoF), center point \mathbf{c} (3-DoF) and a normal vector of the circle $\bar{\mathbf{n}}$ (2-DoF). In the case of the prismatic joint, the direction of the fitting line becomes the direction $\bar{\mathbf{t}}$. In the case of the revolute joint, the axis passes through \mathbf{c} in the direction $\bar{\mathbf{n}}$. Since $\mathbf{l} \cdot \bar{\mathbf{n}} = 0$ holds, the position parameter of the axis \mathbf{l} is;

$$\mathbf{l} = \mathbf{c} - (\mathbf{c} \cdot \bar{\mathbf{n}})\bar{\mathbf{n}}. \quad (4)$$

B. Articulation parameter optimization

The articulation parameters are optimized by aligning the point clouds of the manipulated object through the sequence. Since the point cloud sequence includes both the static background and the dynamic objects, the alignment of the manipulated object needs to ignore the effect of the static background. We assume that points nearby the detected hand are likely to belong to the manipulated object. Therefore, we introduce a weighting scheme on 3D points using the hand position; we set the weighting parameter according to the distance from the hand to each point.

The parameters in the optimization process are the articulation parameters and the amount of movement of object between each frame. For simplicity, let \mathbf{J} be the articulation parameters and m_i be the amount of movement of the manipulated object in i -th frame. In the case of prismatic joint, \mathbf{J} is the moving direction $\bar{\mathbf{t}}$, and m_i is the distance parameter a in Sec. II-B.1. In the case of revolute joint, \mathbf{J} is the rotation axis $\bar{\mathbf{n}}$ and \mathbf{l} , and m_i is the amount of rotation θ in Sec. II-B.2.

We align the point clouds of multiple frames simultaneously by a constrained ICP. Let \mathbf{p}_k be k -th points in i -th frame and \mathbf{q}_k be the closest point of \mathbf{p}_k in j -th frame. Considering constrained ICP alignment with only two frames, i -th and j -th, the error is described as follows:

$$\varepsilon_{i,j} = \sum_k w_k w'_k |((\mathbf{R}_i(\mathbf{J}, m_i)\mathbf{p}_k) + \mathbf{t}_i(\mathbf{J}, m_i)) - ((\mathbf{R}_j(\mathbf{J}, m_j)\mathbf{q}_k) + \mathbf{t}_j(\mathbf{J}, m_j))|, \quad (5)$$

where \mathbf{R}_i is the rotation matrix constrained by the joint parameters \mathbf{J} with the amount of movement of i -th frame m_i , and \mathbf{t}_i is the translation vector as well. w_k, w'_k are the weighting parameters according to the distance from the hand as described above and we define them as follows:

$$w_k = \begin{cases} \left(\frac{1}{C+|\mathbf{p}_k - \mathbf{h}_i|}\right)^2 & (\text{hand position is detected}) \\ 1 & (\text{otherwise}) \end{cases} \quad (6)$$

$$w'_k = \begin{cases} \left(\frac{1}{C+|\mathbf{q}_k - \mathbf{h}_j|}\right)^2 & (\text{hand position is detected}) \\ 1 & (\text{otherwise}) \end{cases} \quad (7)$$

where \mathbf{h}_i is hand position in i -th frame and C is a constant value empirically determined. The optimal parameters $\hat{\mathbf{J}}, (\hat{m}_1 \dots \hat{m}_N)$ are obtained by simultaneously minimizing the error as follows:

$$\hat{\mathbf{J}}, (\hat{m}_1 \dots \hat{m}_N) = \arg \min_{\mathbf{J}, (m_1 \dots m_N)} \sum_{i,j} \varepsilon_{i,j}, \quad (8)$$

where N is the number of frames.

C. Segmentation

We segment the points in the sequence of the point clouds into the manipulated object and the background. After aligning the point clouds as described in the previous section, the areas that move symmetrically about the articulation axis overlap each other, as shown in Fig. 4 (a). We first extract the points in these areas as the potential regions of

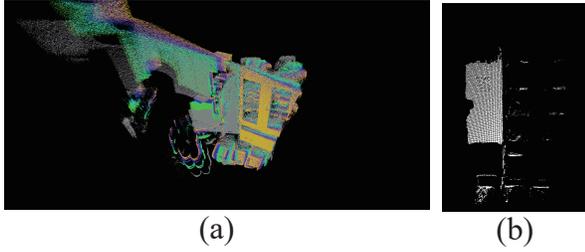


Fig. 4. (a) Point clouds after the alignment of the manipulated object. (b) Extracted symmetric objects.

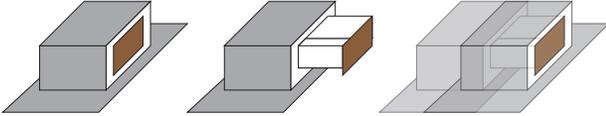


Fig. 5. Prismatic object alignment. After aligning the brown area, gray areas are also detected as a symmetric area because of these overlap between before and after the movement.

the manipulated object, as shown in Fig. 4 (b). Next, since the extracted points include ambiguous regions, we select the regions of the manipulated object by using a clustering method and the hand position by assuming that the object is connected to the hand that manipulates it.

1) *Initial points extraction:* We consider that the points of the manipulated object move symmetrically according to the articulation model through the frames. For example, in the case of the prismatic drawer, the points on the drawer’s surface move in parallel along the direction of movement. In the case of the revolute drawer, such as a door, the points move symmetrically about the axis of rotation. Such areas overlap with each other in the sequential frames after alignment in the previous section. In other words, the error between the corresponding points is small in these areas.

We obtain the points in the overlapping areas using the alignment error. If the distance from a point in each frame to the closest point in the background frame is less than a threshold value, we extract the point as the potential manipulated object’s point. The threshold value depends on the sensor accuracy. In our experiments, after calculating distances from the background frame to each frame, we used the median distance and set the threshold to 5cm for the first time and 3cm for the subsequent iterations in the refinement process described later in Sec. III-D.3.

2) *Region selection:* The extracted point cloud contains confidence and ambiguous regions. The confidence region really moves according to the articulation model. On the other hand, the ambiguous region looks moving according to the model, but it is not easy to identify it to be background or target. For example, in the case of the prismatic joint, a region like a floor in Fig. 5 whose surface normal is perpendicular to \bar{t} is ambiguous. In the case of the revolute

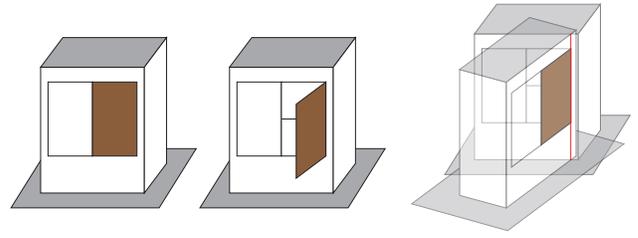


Fig. 6. Revolute object alignment. After aligning the brown area, gray areas are also detected as a symmetric area because of these overlap between before and after the movement.



Fig. 7. Example of symmetric areas of a revolute joint. A part of the floor is detected as a symmetric area due to its symmetric structure on the rotational axis of the chair.

joint, a region like a ceiling or a floor in Fig. 6 that is symmetric about the rotation axis is also ambiguous. Figure 7 shows an example of the symmetric area including ambiguous objects.

Therefore, we need to select points of the manipulated object from the initially extracted point cloud. The confidence region can be identified by the surface normal and the articulation model. The surface normal of the point in the confidence region is the same direction to \bar{t} in the case of the prismatic joint. In the case of the revolute joint, the surface normal moves along the tangential direction about the rotation axis. We use the hand information again to select the points from ambiguous regions by assuming that the manipulated object connects to the hand. We apply Euclidian clustering [21] on the extracted point cloud and select the clusters nearby the hand position in the first frame with a threshold distance.

D. Refinement with a hand soft constraint

We refine the articulation parameters by aligning the segmented point clouds. Since the first alignment process in Sec. III-B uses all points, including non-manipulated objects, the estimated parameters still have room for improvement in the accuracy. On the other hand, in the case of a feature-less object, the segmented point cloud does not have distinctive features for alignment. Therefore, we introduce a soft constraint by the hand in the refinement process.

1) *Geometric error term:* We define the geometric error by the distance of the corresponding points between the

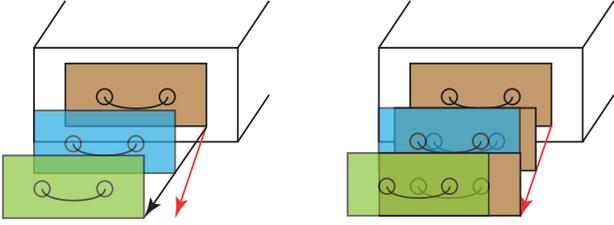


Fig. 8. Alignment ambiguity in a flat prismatic object. Consider the case where the flat drawer moves along with the black arrow, however the estimated direction is the red. The alignment performs with even the wrong estimated direction by treating the non-overlapped parts as outliers of alignment

segmented points in the first frame and the corresponding points in i -th frame. The cost of the geometric error c_{geo} is defined as the same with Eq. 8 with constant weights $w_k = w'_k = 1$ as follows:

$$c_{geo} = \sum_j \varepsilon_{0,j} \quad (9)$$

Note that $\varepsilon_{0,j}$ includes the error of only the segmented points.

2) *Hand soft constraint*: We introduce soft constraint of hand to support the alignment of relatively feature-less surfaces and to eliminate the ambiguous regions. For example, in the case of prismatic objects as shown in Fig. 8, the alignment error in the region where the surface is almost flat in the direction parallel to \vec{t} becomes small even the corresponding points are wrong by the sliding effect.

We use the 2D positions of the hand joints which are also given by the hand-keypoint detector. The soft constraint of hand is the summation of the geometric error of each joint position. The cost of soft constraint of hand c_{hand} is,

$$c_{hand} = \sum_l \sum_j \alpha_{0,l} \alpha_{j,l} |\mathbf{h}_{0,l} - ((\mathbf{R}_j \mathbf{h}_{j,l}) + \mathbf{t}_j)|, \quad (10)$$

where $\mathbf{h}_{j,l}$ is the position of l -th hand joint in j -th frame, and $\alpha_{j,l}$ is its confidence value. The confidence value is also given by the detector.

3) *Optimization*: Finally, the parameters are optimized by minimizing the joint cost described as follows:

$$\hat{\mathbf{J}}, (\hat{m}_1 \dots \hat{m}_N) = \arg \min_{\mathbf{J}, (m_1 \dots m_N)} \frac{1}{n} c_{geo} + \lambda c_{hand}, \quad (11)$$

where n is number of points in the segmented region, and λ is the weight value which is empirically given to adjust the effect of the hand constraint. After the parameter refinement, the segmentation result is also refined by the same process described Sec. III-C.

IV. IMPLEMENTATION

The scanned person and apparent static objects are removed from point clouds to reduce the number of points for computational efficiency and improve the robustness. We apply RCNN-Mask [22] to RGB images to find the region of the person and remove the points in the region. The static

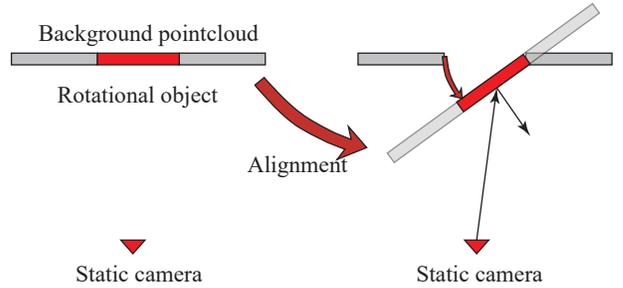


Fig. 9. We remove points in a frame with a shallower angle of incidence from the camera for median computation of distance to the nearest points in each frame.

objects can be removed by comparing the depth values in each pixel of the background frame and the manipulation sequence.

We use a downsampled point cloud in the initial estimation process for computational efficiency. We also limit the number of frames up to 15 frames. If the number of frames is higher than that, we sub-sample 15 frames from the original sequence in which the hand positions are correctly detected. We use kd-tree [23] for searching nearest neighbor points and Voxel Grid Filter to query in the downsampled kd-tree. After the first segmentation, we use all points for the constrained ICP.

In the segmentation process, we also use surface normal information for filtering points out. The accuracy of depth measurement gets worth when an incident angle of light is small. We derive the symmetric region with filtering out the points in a frame with a smaller incident angle from the camera (See Fig 9).

We also have several empirically determined values. We set $C = 0.2$ in Eq. 6 and $\lambda = 0.01$ in Eq. 11. We used the implementations in [24] for the Euclidean clustering for segmentation, Voxel Grid Filter, and kd-tree. We also used ceres-solver [25] for the non-linear optimization.

V. EXPERIMENTAL RESULTS

We first validate the accuracy of the estimated articulation model. We also demonstrate that the proposed method works well in various scenes. We used Microsoft Kinect v2 [26] as the RGB-D input device in all experiments. The length of the original input RGB-D sequences was around 10-40 frames.

A. Accuracy evaluation

We used a flat prismatic object for validating the effect of soft hand constraint. In addition to this, we demonstrate that the proposed method estimates the articulation model and segmentation of a revolute joint object accurately.

1) *Validation of hand soft constraint*: We used a flat drawer and pasted ArUco code on it for the comparative evaluation. Figure 10 shows the estimation result of hand motion fitting (green), ArUco code (yellow), final estimation result with soft hand constraint (red), and without it (light blue). Since the ArUco code detection works well, we can

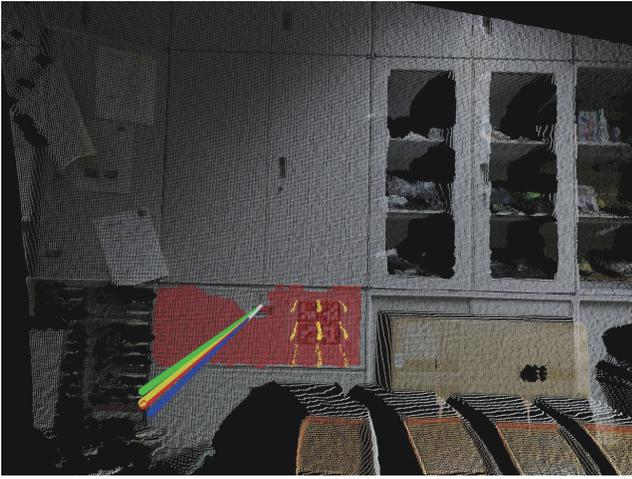


Fig. 10. The articulation and segmentation estimation result with a flat drawer. (Green) Line fitting with the hand positions. (Yellow) Estimation using the trajectory of ArUco code corners, the corner trajectories are also shown as the yellow points. (Blue) The estimation result without hand soft constraint, (Red) The estimation result with hand soft constraint. Red points are the segmentation result with our method.

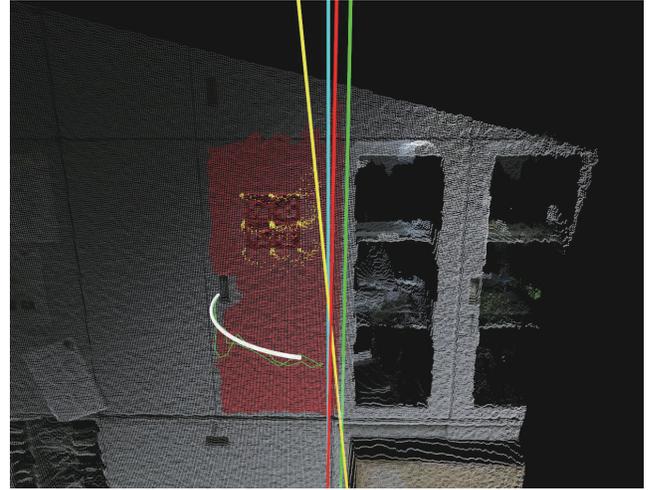


Fig. 11. The articulation and segmentation estimation result with a shelf door. (Light blue) The parameters manually calculated from the hinge structure. (Green) Line fitting with the hand positions. (Yellow) Estimation using the trajectory of ArUco code corners, the corner trajectories are also shown as the yellow points. (Blue) The estimation result without hand soft constraint, (Red) The estimation result with hand soft constraint. Red points are the segmentation result with our method. (White) movement range

assume that the result of the ArUco code is a reference result in this case. Comparing to the ArUco code result, the constrained ICP with hand motion is almost in the same direction, whereas the hand motion fitting and original ICP show an apparent error in the direction. The result indicates that the human hand information works well for the articulation model estimation.

2) *Accuracy evaluation in revolute object:* Figure 11 shows the result of the revolute articulation estimation. In this case, we manually give the reference result by pointing out the hinges (right blue). The estimation result of hand motion fitting (green), ArUco code (yellow) and final segmentation result with hand soft constraint (red) are shown in the figure. The result by ArUco (yellow) has a more significant error due to the distortion when the incident angle becomes small. The circle fitting on the hand motion (green) still has the error in position. After the refinement, the articulation parameters (red) have been clearly improved.

B. Other results

Figure 12 shows the results of several scenes, including both prismatic or revolute objects. Only for the prismatic object in the top-right in Fig. 12, we applied the hand region detector [19] since there was not enough body area in the picture for detecting hand positions by OpenPose. One of the advantages of the proposed method is robustness. In the initial estimation by the hand motion, the proposed method correctly identified the articulation type of the objects. In all scenes of Fig. 12, the final segmentation results demonstrate that our method correctly estimated the manipulated object's region without ambiguous regions.

VI. CONCLUSION

In this paper, we presented a hand-motion-guided articulation model estimation and segmentation method. The proposed method uses the hand position information for the initial articulation estimation, the model refinement by a constrained ICP alignment, and selecting the region of the manipulated object from ambiguous regions. The experimental results demonstrate that our method correctly estimates the articulation model and segmentation of manipulated objects. For future work, we will extend the method and apply it for articulated objects with more degrees of freedom.

REFERENCES

- [1] D. A. Ross, D. Tarlow, and R. S. Zemel, "Learning articulated structure and motion," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 214–237, 2010.
- [2] B. Jacquet, R. Angst, and M. Pollefeys, "Articulated and restricted motion subspaces and their signatures," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2013.
- [3] K. Pauwels, L. Rubio, and E. Ros, "Real-time model-based articulated object pose detection and tracking with variable rigidity constraints," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [4] X. Huang, I. Walker, and S. Birchfield, "Occlusion-aware reconstruction and manipulation of 3d articulated objects," in *2012 IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 1365–1371.
- [5] S. Pillai, M. R. Walter, and S. Teller, "Learning articulated motions from visual demonstrations," in *Proceedings of Robotics: Science and Systems (RSS)*, Berkeley, CA, July 2014.
- [6] Y. Pekelny and C. Gotsman, "Articulated object reconstruction and markerless motion capture from depth video," in *Computer Graphics Forum*, vol. 27, no. 2. Wiley Online Library, 2008, pp. 399–408.
- [7] D. Tzionas and J. Gall, "Reconstructing articulated rigged models from rgb-d videos," in *European Conference on Computer Vision*. Springer, 2016, pp. 620–633.

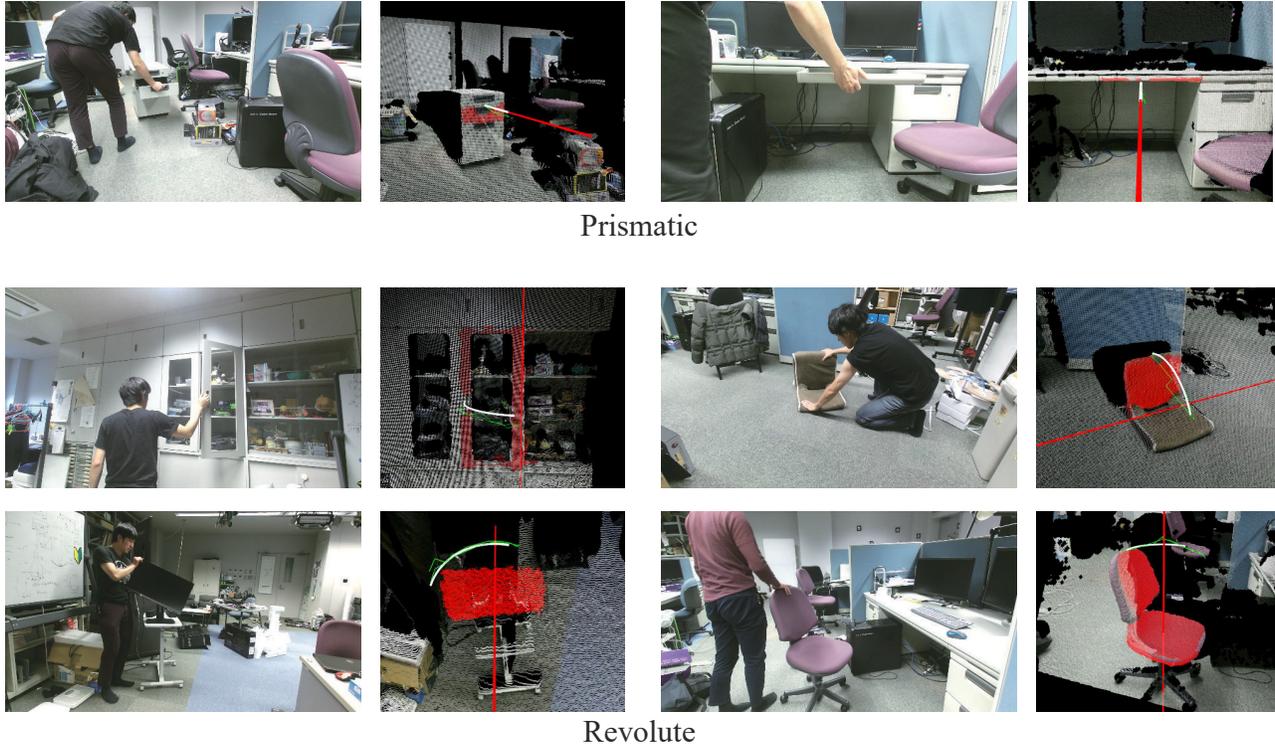


Fig. 12. Other manipulation scenes and articulation and segmentation results. The articulation estimation results are shown as the red line. The segmentation results also shown as the red parts. Hand trajectory is also shown as the green polygonal line. The white lines indicate movement range in the video sequence. The top two are prismatic joint and the rest are revolute joint.

- [8] K. Chaudhary, K. Okada, M. Inaba, and X. Chen, "Predicting part affordances of objects using two-stream fully convolutional network with multimodal inputs," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018.
- [9] A. F. Daniele, T. M. Howard, and M. R. Walter, "A multiview approach to learning articulated motion models," in *Robotics Research*. Springer, 2020, pp. 371–386.
- [10] K. Hausman, S. Niekum, S. Osentoski, and G. S. Sukhatme, "Active articulation model estimation through interactive perception," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2015, pp. 3305–3312.
- [11] J. Sturm, K. Konolige, C. Stachniss, and W. Burgard, "Vision-based detection for learning articulation models of cabinet doors and drawers in household environments," in *2010 IEEE International Conference on Robotics and Automation*. IEEE, 2010, pp. 362–368.
- [12] J. Fayad, C. Russell, and L. Agapito, "Automated articulated structure and 3d shape recovery from point correspondences," in *2011 International Conference on Computer Vision*. IEEE, 2011, pp. 431–438.
- [13] Y. Kim, H. Lim, S. C. Ahn, and A. Kim, "Simultaneous segmentation, estimation and analysis of articulated motion from dense point cloud sequence," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 1085–1092.
- [14] X. Lu, H. Chen, S.-K. Yeung, Z. Deng, and W. Chen, "Unsupervised articulated skeleton extraction from point set sequences captured by a single depth camera," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [15] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh, "OpenPose: realtime multi-person 2D pose estimation using Part Affinity Fields," in *arXiv preprint arXiv:1812.08008*, 2018.
- [16] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [17] M. Schroder and H. Ritter, "Hand-object interaction detection with fully convolutional networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017.
- [18] T. Simon, H. Joo, I. Matthews, and Y. Sheikh, "Hand keypoint detection in single images using multiview bootstrapping," in *CVPR*, 2017.
- [19] D. Victor, "Handtrack: A library for prototyping real-time hand tracking interfaces using convolutional neural networks," *GitHub repository*, 2017. [Online]. Available: <https://github.com/victordibia/handtracking/tree/master/docs/handtrack.pdf>
- [20] S. Bambach, S. Lee, D. J. Crandall, and C. Yu, "Lending a hand: Detecting hands and recognizing activities in complex egocentric interactions," in *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [21] R. B. Rusu, "Semantic 3d object maps for everyday manipulation in human living environments," *KI-Künstliche Intelligenz*, vol. 24, no. 4, pp. 345–348, 2010.
- [22] W. Abdulla, "Mask r-cnn for object detection and instance segmentation on keras and tensorflow," <https://github.com/matterport/Mask-RCNN>, 2017.
- [23] D. M. Mount and S. Arya, "Ann: A library for approximate nearest neighbor searching," <https://www.cs.umd.edu/~mount/ANN/>.
- [24] R. Rusu and S. Cousins, "3d is here: Point cloud library (pcl)," in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, May 2011, pp. 1–4.
- [25] S. Agarwal, K. Mierle, and Others, "Ceres solver," <http://ceres-solver.org>.
- [26] "Microsoft Kinect for Xbox One," <http://www.xbox.com/en-US/xbox-one/accessories/kinect-for-xbox-one>.