# The RBO Dataset of Articulated Objects and Interactions

**Roberto Martín-Martín\*, Clemens Eppner\* and Oliver Brock**

## Abstract

We present a dataset with models of 14 articulated objects commonly found in human environments and with RGB-D video sequences and wrenches recorded of human interactions with them. The 358 interaction sequences total 67 minutes of human manipulation under varying experimental conditions (type of interaction, lighting, perspective, and background). Each interaction with an object is annotated with the ground truth poses of its rigid parts and the kinematic state obtained by a motion capture system. For a subset of 78 sequences (25 minutes), we also measured the interaction wrenches. The object models contain textured three-dimensional triangle meshes of each link and their motion constraints. We provide Python scripts to dow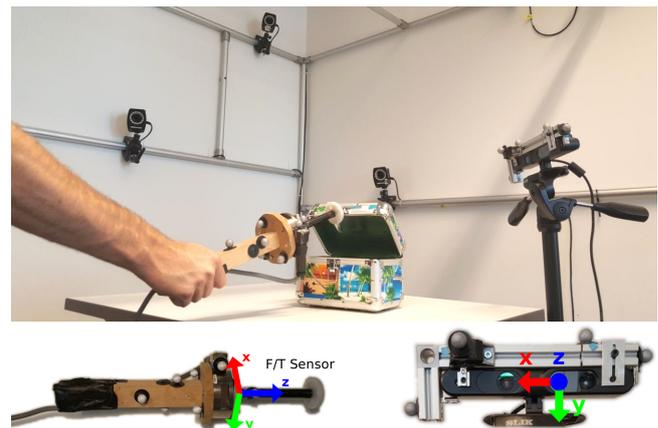nload and visualize the data. The data is available at https://tu-rbo.github.io/articulated-objects/ and hosted at https://zenodo.org/record/1036660/.

arXiv:1806.06465v1 [cs.RO] 17 Jun 2018

## Introduction

The RBO dataset is a collection of 358 RGB-D video sequences (67 minutes) of humans manipulating 14 articulated objects under varying exeperimental conditions (type of interaction, lighting, perspective, and background). All sequences are annotated with ground truth of the poses of the rigid parts and the kinematic state of the articulated object (joint states) obtained with a motion capture system. We also provide kinematic models of these objects including three-dimensional textured shape models. For 78 sequences (25 minutes) the interaction wrenches during the manipulation are also recorded.

We present the first dataset with *articulated* objects. All similar datasets contain *single rigid-body* objects that move or are being manipulated. There are two datasets that could be considered close to ours. The first one (Garcia Cifuentes et al. 2017) is a dataset that was released together with a method for robot arm tracking. This dataset contains images and joint encoder values of a moving robot arm and its kinematic and shape model. In contrast, our dataset is targeted to the study of interactions with everyday human objects: it contains models of multiple ubiquitous articulated objects and sequences of interactions in varying environmental conditions. The second dataset (Michel et al. 2015) provides models of four everyday articulated objects and ground truth of the static pose of each link and the changing pose of camera during the video sequence. The sequences of this dataset do not contain any interaction or manipulation of the objects, only camera motion. Therefore, neither of these datasets can be used to study interactions and to evaluate methods for perceiving them.

Our dataset will help evaluate algorithms for tracking articulated objects (e.g. Schmidt et al. (2014)) and building models of articulated objects (e.g. Martín-Martín et al. (2016)). Apart from benchmarking, the dataset can also be



**Figure 1.** Our sensor setup for recording interactions with articulated objects; Top: interaction in the motion capture volume; Bottom left: tool with F/T sensor for recording interaction wrenches, motion capture markers, and measurements reference frame; Bottom right: Asus RGB-D sensor with motion capture markers and measurements reference frame

used to develop data-driven algorithms by exploiting the provided models to generating virtual visual data. While there is a vast amount of three-dimensional models and datasets of objects (Kasper et al. 2012; Calli et al. 2015), very few of them include and describe articulated mechanisms.

All authors are with the Robotics and Biology Laboratory, Technische Universität Berlin, Germany. \* denotes equal contribution.

**Corresponding author:**
Roberto Martín-Martín, Robotics and Biology Laboratory, Marchstr. 23, 10587 Berlin, Germany
Email: roberto.martinmartin@tu-berlin.de

## Sensor Setup

We use the following sensors to record human interactions with articulated objects (see Fig. 1):

- RGB-D camera Asus Xtion Pro Live, 640×480 pixels, 30 Hz, pointed at the object.
- Motion capture system by Motion Analysis (2017), capture volume: $1.5\,m^3$, including 18 Osprey cameras, providing 3-D positions of fiducial markers at 100 Hz.
- F/T sensor ATI FTN-Gamma DAQ/Net, calibration SI-130-10, force/torque resolution: $F_x = F_y = \frac{1}{40}N$, $F_z = \frac{1}{20}N$, $T_x = T_y = T_z = \frac{1}{800}Nm$, recorded at 100 Hz.

We acquired separately 3-D triangle meshes of each articulated object part, using the following sensors and methods:

- Structured light scanning system by David (2017), SLS-3 HD, scan size: 60-500 mm, resolution: down to 0.05 mm.
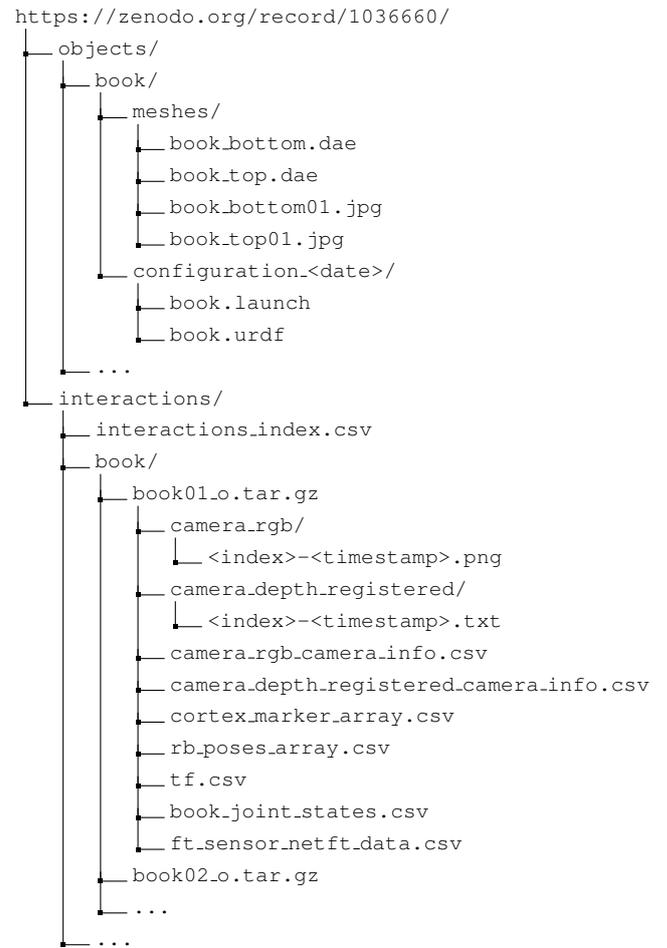- Reconstruction software Autodesk ReMake (2017), used with photos taken with Casio Exilim EX-FC100, 9 MP.

## Data Structure and Usage

The RBO dataset is available at https://tu-rbo.github.io/articulated-objects/ and hosted at https://zenodo.org/record/1036660/. It is composed of two parts (see Fig. 2): a first part with descriptions of 14 articulated objects and the main part containing 358 human interactions with these objects. We also provide Python scripts to facilitate downloading, visualizing and working with the data.

### *Models of Articulated Object (`objects/`)*

The RBO dataset contains 14 models of articulated mechanisms that are commonly found in human environments. Table 1 depicts these objects and their kinematic structure. Each object model (`<object_id>/`) consists of:

- **Link geometries** (`meshes/`): We describe the shape of a link as a three-dimensional triangle textured mesh in the COLLADA (2012) format (`<part_name>.dae`). We provide the associated texture as JPEG images (`<part_name><index>.jpg`).
- **Kinematic structure** (`configuration_<date>/`): We define the relation of links and joints with the widely-used *Unified Robot Description Format* URDF (2017) (`<object_id>.urdf`), an XML file format to describe all elements of articulated objects with chain or tree structure. The objects of the database possess one degree-of-freedom (DoF) joints that can be either prismatic or revolute joints. The base link is the origin of the kinematic tree or chain. It is either rigidly connected to the environment (represented with a *static* joint with zero DoF) or completely unconstrained (represented with a floating joint with six DoF). The reference coordinate frame of a link corresponds to a marker set of the motion capture system (see Section *Data Acquisition*). We define

```
https://zenodo.org/record/1036660/
   objects/
      book/
         meshes/
            book_bottom.dae
            book_top.dae
            book_bottom01.jpg
            book_top01.jpg
         configuration_<date>/
            book.launch
            book.urdf
      ...
   interactions/
      interactions_index.csv
      book/
         book01_o.tar.gz
            camera_rgb/
               <index>-<timestamp>.png
            camera_depth_registered/
               <index>-<timestamp>.txt
            camera_rgb_camera_info.csv
            camera_depth_registered_camera_info.csv
            cortex_marker_array.csv
            rb_poses_array.csv
            tf.csv
            book_joint_states.csv
            ft_sensor_netft_data.csv
         book02_o.tar.gz
         ...
      ...
```

**Figure 2.** The dataset is structured by objects and interactions. Please see text for details.

the joint parameters and link meshes with respect to these coordinate frames. Since marker set locations can vary between recording sessions, we provide a separate kinematic structure description for each session indicated by the `_<date>` suffix in the folder name.

### *Interactions (`interactions/`)*

The RBO dataset contains sensor data of 358 human interactions with the 14 modelled objects ($\geq 25$ interactions per object). The sequences last between $2.7\,s$ and $69.0\,s$ (median: $9.15\,s$). They differ in lighting conditions, camera perspective and motion, background, clutter, actuation of the mechanisms and human motion (see Table 2). The file `interactions_index.csv` contains a list of all interactions and their properties.

The sensor data is organized per object (`<object_id>/`) and interaction (`<object_id><index>_o/`). Each interaction includes:

- **RGB images**: We store the color images as 8-bit loss-less compressed PNG files (`camera_rgb/<index>-<timestamp>.png`).
- **Depth images**: We register the depth images to the RGB camera frame (see Section *Data Acquisition*) and store them as text files containing distances in meters (`camera_depth_registered/<index>-<timestamp>.txt`).

**Table 1.** The dataset contains 14 articulated objects. Letters in the last column represent (S)tatic, (F)loating, (R)evolute and (P)rismatic joints.

| Object ID | Picture | Actuated Joints |
|---|---|---|
| globe | | F — R |
| laptop | | F — R |
| ikea | | S < R / P |
| foldingrule | | F — R — R |
| book | | F — R |
| treasurebox | | F — R |
| tripod | | F — P — R |
| clamp | | F — P |
| pliers | | F — R |
| cardboardbox | | F — R |
| rubikscube | | F — R |
| microwave | | S — R |
| ikeasmall | | S < P / P |
| cabinet | | S < P / P |

**Table 2.** Properties of the 358 interaction recordings

| | | |
|---|---|---|
| Lighting Conditions | Artificial | 178 |
| | Natural | 107 |
| | Dark | 73 |
| Camera Motion | Yes | 100 |
| | No | 258 |
| Type of Background | Plain | 197 |
| | Textured | 91 |
| | Black | 70 |
| Clutter | Yes | 109 |
| | No | 249 |
| Actuated DoFs | Only internal | 162 |
| | Internal and external | 196 |
| Interaction Wrenches | Yes | 78 |
| | No | 280 |

- **Intrinsic camera parameters**: We provide the focal length, center point and distortion parameters of the *Plumb Bob* model (Brown 1966) for both, the camera that generates RGB (`camera_rgb_camera_info.csv`) and the one that generates depth images (`camera_depth_registered_camera_info.csv`).
- **Extrinsic camera parameters**: We represent the 6-D transformation between the cameras of the RGB-D sensor with a translation vector and a quaternion (`tf.csv`).
- **Infra-red marker positions**: We include the 3-D locations of all infrared fiducial markers in the scene measured by the motion capture system (`cortex_markers_array.csv`).
- **Rigid body poses**: We store the 6-D poses defined by sets of infrared fiducial markers as position vectors and quaternions. We include the pose of each link of the articulated object, the RGB-D and force/torque (F/T) sensor at 100 Hz (`rb_poses_array.csv`).
- **Joint configurations**: We compute the joint configuration of the articulated object in the scene from the 6-D poses of its links (`<object>_joint_states.csv`).
- **Wrenches**: We provide the forces and torques for the interaction as provided by the F/T sensor (`ft_ssensor_netft_data.csv`). We include this data in at least five interactions per object.

## Utilities

We provide Python scripts and a ROS package on the website https://tu-rbo.github.io/articulated-objects/ to facilitate the download and visualization of the data.

- The download script (`rbo_downloader.py`) fetches object models and interaction files. The user can also select groups of interactions fulfilling a certain property, e.g. all interactions with an object, or all interactions with wrench measurements.
- The visualization script (`rbo_visualizer.py`) displays the content of an interaction folder: RGB, depth images, wrenches and/or joint states.
- ROS package: Additionally to the interaction sequences in the file format described above, we also provide all data in the form of a ROSBags (2017). We

**Figure 3.** Visualization of the sensor data for opening a drawer: RGB and depth image (*left*), mesh model, point cloud and coordinate frames of tracked bodies (*middle*), and drawer state and applied wrenches (*right*) with dashed vertical line indicating current time.

provide a ROS package including scripts to visualize the data in this format.

## Data Acquisition

### Visual Data and 6-D Body Poses

The main goal of our dataset is to evaluate and develop algorithms based on visual data (RGB or RGB-D) with/without interaction wrenches for the perception of articulated objects. For this goal, it is crucial to register accurately the visual information and the ground truth provided by the motion capture system. We first calibrate the intrinsic parameters of the RGB-D sensor. We use a checkerboard of known dimensions and take pictures at different poses of the checkerboard with respect to the camera with both the RGB and the infrared camera of the RGB-D sensor. We use an OpenCV-based camera calibration tool to estimate the internal parameters of the cameras (focal length, center point, and distortion parameters of the Plumb model) by detecting corner points on the checkerboard and estimating the parameters that minimize the squared error of the reprojection of the points from a separate multi-view PnP procedure per camera. We then rectify the color and infrared images and estimate the 6-dimensional (6-D) transformation between the RGB and the infrared camera of the RGB-D sensor from a multi-view PnP procedure between the cameras.
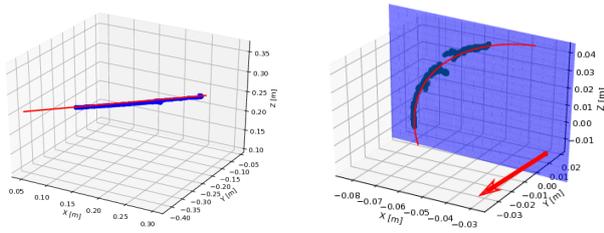
We calibrate the extrinsic parameters of the RGB-D sensor with respect to the set of motion capture markers attached to it (see Fig. 1). We attached infrared markers to the corners of the checkerboard and use the motion capture system to detect their 3-D location. The points are projected on the color image based on the currently estimated transformation between the sensor and the marker set. We minimize the error of the projection of the point markers on the color image at different locations of the checkerboard. After the calibration procedure, the point clouds recorded from the RGB-D sensor are registered with respect to the motion capture readings.

To acquire 6-D pose measurements of the articulated object from the motion capture system we attach marker sets to each of the links and place them inside the tracking volume. The motion capture system estimates the 3-D location of the infrared fiducial markers on the scene with submillimeter accuracy and generates a 6-D pose measurement based on the predefined model of the arrangement of the markers within each marker set. While a minimum of three infrared markers defines a marker set we use at least five markers per set to improve the accuracy of the 6-D pose measurements and the robustness of the system against occlusions.

### Kinematic Properties

We use the 6-D rigid body poses to compute the ground truth of the joint parameters in an offline batch procedure. We estimate the axis of a prismatic joint by fitting a line to the time-varying positions of the child body with respect to the parent body in a least-squares sense (Fig. 4, left). The computation of the kinematic state of a prismatic joint requires to define an origin. Without loss of generality we use the first pose of the child body with respect to the parent body as the origin. We calculate the prismatic joint state as the distance between the child body and this point along the fitted line.

For revolute joints, we estimate the orientation of the joint axis as a unit vector, and its position as a 3-D point. We obtain the plane that best fits the positions of the child body with respect to the parent body during the interaction by minimizing the squared distance of the points to the plane. The plane's normal corresponds to the orientation of the revolute axis (Fig. 4, right). We then project all the points to the plane and estimate a circle using a least-squares fit with respect to the projected points. The circle's center indicates the position of the revolute axis. Without loss of generality we use the radius connecting the projection of the first pose of the child body with respect to the parent body as origin. We calculate the configuration of the revolute joint as the

**Figure 4.** Estimating the joint axis for a prismatic (*left*) and revolute joint (*right*): The blue dots show the position of the moving body, the red lines are the fitted axes. For the revolute joint (*right*) the orientation of the axis is obtained by fitting a plane (blue), while its position is based on a circle fit (red) within that plane.

angle of the arc between the child body and this point along the circle defined by joint axis position and orientation.

## Interaction Wrenches

For each object we provide five interactions with measurements of the interaction wrench. In these interactions the humans actuate the articulated object with a tool attached to a force/torque (F/T) sensor (Figure 1). The motion capture system measures the pose of the interaction tool and we provide it as part of the interaction data. We also provide a three-dimensional textured model of the tool with the sensor. The wrenches measured by the sensor and included in the dataset are raw values with bias. We measured the bias in the measurements with a calibration procedure where we align in turns one of the main axis of the F/T sensor to the vertical direction and collect the sensor readings. The result of the calibration is the following wrench bias vector:

$$w_b = \begin{pmatrix} f_b \\ \tau_b \end{pmatrix} = \begin{pmatrix} (-0.927\,\mathrm{N}, 1.122\,\mathrm{N}, 1.332\,\mathrm{N}) \\ (0.104\,\mathrm{N\,m}, 0.027\,\mathrm{N\,m}, -0.033\,\mathrm{N\,m}) \end{pmatrix}$$

In order to subtract the effect of the tool from the wrench readings we measured the following dynamic properties of the elements attached to the F/T sensor

- Mass: 163.0 g
- Center of Mass: $(0\,\mathrm{cm}, 0\,\mathrm{cm}, -1.5\,\mathrm{cm})$

and the transformation between the frame tracked by the motion capture system and the measurements frame depicted in Fig. 1:

$$T^{ft\_meas}_{ft\_mocap} = \begin{pmatrix} 0.991 & 0.052 & -0.123 & 0.019 \\ -0.060 & -0.650 & -0.757 & -0.008 \\ -0.120 & 0.758 & -0.642 & -0.001 \\ 0.0 & 0.0 & 0.0 & 1.0 \end{pmatrix}$$

## Link Geometries

We use three alternative methods to generate three-dimensional triangle meshes of the links of the articulated models. For small-scale objects, we use a system based on structured light (David 2017). It projects a known light pattern onto the object to generate 3-D information. The scanner acquires partial 3-D models from 24 different view points using a rotating plate and integrates them into a colored triangle mesh. Before the scan we perform an initial calibration procedure to segment background from

foreground. For large-scale textured objects we use the software Autodesk ReMake (2017), which reconstructs a high definition 3-D mesh by applying a multi-view geometric algorithm on overlapping color photos of the object. We take $\approx 25$ photos per object with a Casio Exilim EX-FC100 (resolution: 9 megapixel), located on a hemisphere centered around the object. For large-scale textureless objects like the cabinet we generate meshes by hand using the 3-D creation suite Blender. We post-process all models to fill holes, to remove parts of the surrounding environment, and to register them to the attached motion capture markers.

## References

Autodesk ReMake (2017) Autodesk remake. http://remake.autodesk.com/about. Accessed: 2017-05-30.

Brown DC (1966) Decentering distortion of lenses. *Photogrammetric Engineering and Remote Sensing* .

Calli B, Singh A, Walsman A, Srinivasa S, Abbeel P and Dollar AM (2015) The YCB object and model set: Towards common benchmarks for manipulation research. In: *Proceedings of the 2015 International Conference on Advanced Robotics (ICAR)*. IEEE, pp. 510–517.

COLLADA (2012) Industrial automation systems and integration – COLLADA digital asset schema specification for 3d visualization of industrial data. ISO ISO/PAS 17506:2012, International Organization for Standardization, Geneva, Switzerland.

David (2017) SLS-3 HD. http://hp.com/go/3dscan. Accessed: 2017-05-30.

Garcia Cifuentes C, Issac J, Wüthrich M, Schaal S and Bohg J (2017) Probabilistic articulated real-time tracking for robot manipulation. *IEEE Robotics and Automation Letters (RA-L)* 2(2): 577–584.

Kasper A, Xue Z and Dillmann R (2012) The KIT object models database: An object model database for object recognition, localization and manipulation in service robotics. *The International Journal of Robotics Research (IJRR)* 31(8): 927–934.

Martín-Martín R, Höfer S and Brock O (2016) An Integrated Approach to Visual Perception of Articulated Objects. In: *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*. Stockholm, Sweden: IEEE, pp. 862–869.

Michel F, Krull A, Brachmann E, Yang MY, Gumhold S and Rother C (2015) Pose estimation of kinematic chain instances via object coordinate regression. In: *Proceedings of the 2015 British Machine Vision Conference (BMVC)*. pp. 181–1.

Motion Analysis (2017) Motion analysis corporation. http://ftp.motionanalysis.com. Accessed: 2017-05-30.

ROSBags (2017) Robot operating system (ros), bags of messages. http://wiki.ros.org/Bags. Accessed: 2017-05-30.

Schmidt T, Newcombe RA and Fox D (2014) DART: Dense articulated real-time tracking. In: *Robotics: Science and Systems*. Berkeley, California, USA, pp. 342–350.

URDF (2017) Robot operating system (ROS), unified robot description format (URDF). http://wiki.ros.org/urdf. Accessed: 2017-05-30.