

# Supplementary Material: A Pipeline for Generating Ground Truth Labels for Real RGBD Data of Cluttered Scenes

Anonymous Author(s)

Affiliation

Address

email

## 1 Comparison with Human Labeling of Single Frame

To approximately quantify the quality of the data generated by our pipeline, and the speed of labeling, we compared with a traditional technique of labeling one image with a polygon of the segmented object (Figure 1). We randomly chose two images from our dataset, and used [1] to label them by hand. A side-by-side comparison of the human labeling and the label generated from our pipeline are provided below. Human labeling using [1] took approximately 10 minutes per frame. With our method we spent approximately 60 seconds of human input to label each of these scenes, but this is amortized over 1,000 views of this scene. Accordingly the human time per label (Figure 1, bottom row) is approximately four orders of magnitude less for our method.

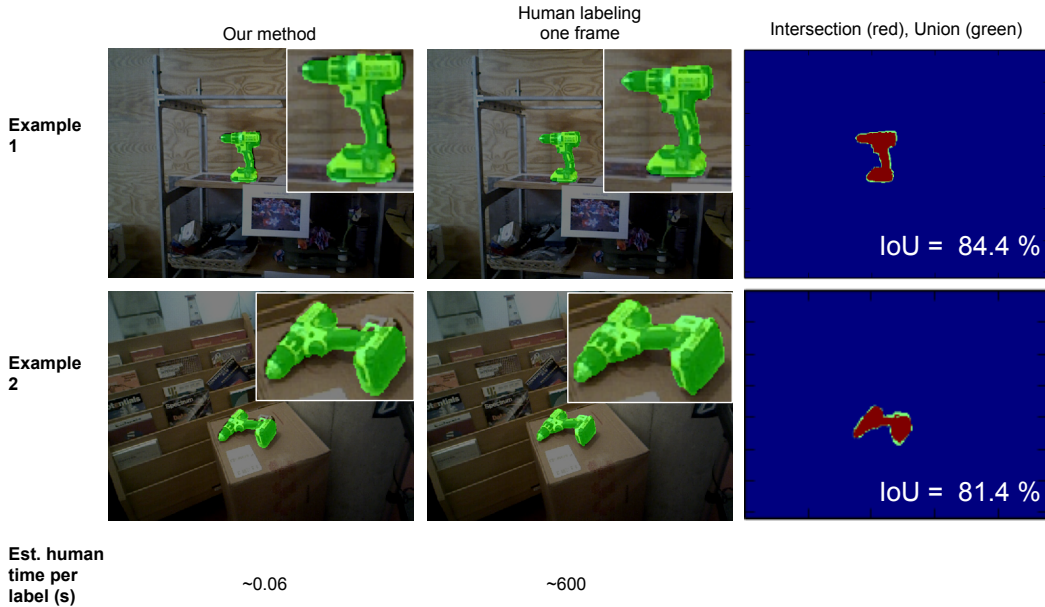


Figure 1: Comparisons of our method (left) vs. human single-frame labeling (middle), for two example images. The insets show a zoomed-in cropped view around the object. For the images, we respectively compute the intersection over union (IoU) for the drill mask at 84.4% and 81.4% (right). Our method has approximately a four orders of magnitude advantage in terms of human labeling time per frame (bottom row).

## 10 2 Experimental Details

11 Here we expand on experimental details that due to space constraints we could not fit into the main  
12 manuscript.

### 13 2.1 Object Set and Object Meshes

14 In total we used a set of 12 object meshes. We tested a variety of methods for acquiring object  
15 meshes. The highest quality meshes we produced (oil bottle, phone, red robot, and drill) were from  
16 our handheld Artec 3D scanner. We also used a tabletop, spinning 3D scanner (toothpaste) which  
17 was more difficult to use. For the tissue box, we simply measured it by hand and created a perfect  
18 box model mesh. Other meshes were obtained from others' datasets, including the blue funnel from  
19 [2] and the cracker box, tomato soup, spam, and mug from the YCB object set [3].

### 20 2.2 Segmentation Network Training

21 Used a TensorFlow reimplementation [4] of DeepLab [5], but without the CRF post-processing step.  
22 We implemented CRF post-processing but found this to not improve results, due to the challenging  
23 occlusions and neighboring objects with similar color textures. We began training our models with  
24 the weights provided by [4], which were pre-trained on the PASCAL VOC dataset. All images from  
25 the native Asus Xtion resolution,  $640 \times 480$ , were downsized to  $480 \times 360$  for training. Parameters  
26 used for training were the defaults in [4] for full-network training:  $2.5 \times 10^{-4}$  step size, 0.9 momentum,  
27 20,000 steps, 2 batch size. The one exception was our "how many views?" experiment, for which  
28 100,000 steps at a batch size of 2 was used in order to allow the potential benefit of the larger datasets.  
29 All models were trained on a GTX 1080; training completed in approximately 2.5 hours for 20,000  
30 steps, and 12.5 hours for 100,000 steps.

### 31 2.3 Training Data Details

32 The empirical evaluations were performed with variations on two primary groups of training data.  
33 These two groups of training data did not encompass all of our data we generated (for example, they  
34 did not include all of the objects we have generated data for, or all the environments), but they did  
35 comprise a focused subset which allowed focused comparisons.

#### 36 2.3.1 Multi-object Training Set

37 A set of 51 total scenes were used for the multi-object training experiments. All of these scenes  
38 were taken the same day, over the course of a few hours, and were each taken with the same pre-  
39 programmed motion of the Kuka IIWA arm with mounted Asus Xtion camera. Lighting was kept  
40 constant throughout all experiments. From each scene were taken  $135 \pm 1$  seconds of data at 30 Hz,  
41 giving 4,000 frames per scene. The six objects for these experiments were: oil bottle, drill, tissue  
42 box, spam, cracker box, and the blue funnel.

43 In summary, the total number of training and test scenes available were:

- 44 • 18 single-object training scenes (3 scenes each for each of 6 objects)
- 45 • 18 multi-object training scenes (each with all 6 objects)
- 46 • 6 single-object test scenes (1 scene each for each of 6 objects)
- 47 • 9 multi-object test scenes

48 More visuals of these scenes are provided in our supplementary video.

#### 49 2.3.2 Drill Training Set

50 A set of 61 total scenes were used for the experimentation with a wide variety of backgrounds.  
51 These scenes were mostly taken by handheld data collection, except for 3 that were from an KUKA-  
52 arm-mounted data collection. Each handheld scene comprised of 34 seconds of 30 Hz data for  
53 approximately 1,000 frames.

54 More visuals of these scenes are provided in our supplementary video.

## 55 References

- 56 [1] P. Tangseng, Z. Wu, and K. Yamaguchi. Looking at outfit to parse clothing. Mar 2017. URL  
57 <http://arxiv.org/abs/1703.01386v1>.
- 58 [2] J. M. Wong, V. Kee, T. Le, S. Wagner, G.-L. Mariottini, A. Schneider, L. Hamilton, R. Chipalkatty,  
59 M. Hebert, D. M. S. Johnson, J. Wu, B. Zhou, and A. Torralba. SegICP: Integrated Deep Semantic  
60 Segmentation and Pose Estimation. *ArXiv e-prints*, Mar. 2017.
- 61 [3] B. Calli, A. Walsman, A. Singh, S. Srinivasa, P. Abbeel, and A. M. Dollar. Benchmarking in  
62 manipulation research: Using the yale-cmu-berkeley object and model set. *IEEE Robotics &*  
63 *Automation Magazine*, 22(3):36–52, 2015.
- 64 [4] DrSleep. DeepLab-ResNet-TensorFlow, [https://github.com/DrSleep/](https://github.com/DrSleep/tensorflow-deeplab-resnet)  
65 [tensorflow-deeplab-resnet](https://github.com/DrSleep/tensorflow-deeplab-resnet). [https://github.com/DrSleep/](https://github.com/DrSleep/tensorflow-deeplab-resnet)  
66 [tensorflow-deeplab-resnet](https://github.com/DrSleep/tensorflow-deeplab-resnet).
- 67 [5] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic  
68 image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs.  
69 *arXiv:1606.00915*, 2016.