

# Appendix of ImVoxelNet

Danila Rukhovich<sup>1,2</sup>, Anna Vorontsova<sup>1</sup>, Anton Konushin<sup>1,2</sup>

<sup>1</sup>Samsung AI Center Moscow; <sup>2</sup>Lomonosov Moscow State University

{d.rukhovich, a.vorontsova, a.konushin}@samsung.com

## Abstract

*In A, we report the results of the extensive evaluation of the proposed method on the SUN RGB-D dataset. In Section B we visualize predicted bounding boxes for several samples taken from all four datasets we use in our experiments.*

encoding is consistent within each dataset.

## A. More results on SUN RGB-D

For a comprehensive comparison, we also mention PerspectiveNet [4], which is evaluated following a different protocol. In that protocol, the annotations are mapped into 30 object categories. Accordingly, we train ImVoxelNet using the same object categories. The results are reported in Tab. 3. Among these 30 categories, 10 object categories are consistent with 10 categories used in [3, 2, 5]. So, we can merge these benchmarks and report metrics for [3, 2, 5, 4] that are obtained on the same subset of 10 object categories. Following [4], we assume camera poses are known, so we optimize only  $L_{indoor}$  and do not use any additional camera pose loss.

Another SUN RGB-D benchmark has been proposed in [7] for point cloud-based methods evaluation. This benchmark implies detecting objects of 10 categories with mAP@0.25 chosen as the main metric. In Tab. 1, we report the results of our method against point cloud-based methods. This comparison is unfair, favoring point cloud-based methods since they have access to more complete data. Nevertheless, we report the metrics to establish a baseline for monocular 3D object detection on SUN RGB-D.

Comparison with Total3DUnderstanding [5] on all NYU-37 object categories is present in Tab. 4. In this experiment, we optimize  $L_{indoor} + L_{extra}$  since camera pose is assumed unknown.

## B. Visualization

All visualized images belong to validation subsets of the corresponding datasets. Different colors of the depicted bounding boxes mark different object categories; the color

Method	RGB	PC	bath	bed	bkshf	chair	desk	dresser	nstand	sofa	table	toilet	mAP
F-PointNet[6]	✓	✓	43.3	81.1	33.3	64.2	24.7	32.0	58.1	61.1	51.1	<b>90.9</b>	54.0
VoteNet[7]	✗	✓	74.4	83.0	28.8	75.3	22.0	29.8	62.2	64.0	47.3	90.1	57.7
H3DNet[9]	✗	✓	73.8	85.6	31.0	<b>76.7</b>	<b>29.6</b>	33.4	65.5	66.5	50.8	88.2	60.1
ImVoteNet[8]	✓	✓	<b>75.9</b>	<b>87.6</b>	<b>41.3</b>	<b>76.7</b>	28.7	<b>41.4</b>	<b>69.9</b>	<b>70.7</b>	<b>51.1</b>	90.5	<b>63.4</b>
ImVoxelNet	✓	✗	71.7	69.6	5.7	53.7	21.9	21.2	34.6	51.5	39.1	76.8	40.7

Table 1. AP@0.25 scores for 10 object categories [7] from the SUN RGB-D dataset. All methods but ImVoxelNet use point cloud (PC) as an input.

Method	bed	chair	sofa	table	desk	toilet	bin	sink	shelf	lamp	mAP
3DGP[1]	5.62	2.31	3.24	1.23	-	-	-	-	-	-	-
HoPR[3]	58.29	13.56	28.37	12.12	4.79	16.50	0.63	2.18	1.29	2.41	14.01
CooP[2]	63.58	17.12	41.22	26.21	9.55	58.55	10.19	5.34	3.01	1.75	23.65
PerspectiveNet[4]	<b>79.69</b>	40.42	62.35	44.12	20.19	81.22	22.42	<b>41.35</b>	8.29	13.14	39.09
ImVoxelNet	77.87	<b>65.94</b>	<b>63.89</b>	<b>51.17</b>	<b>31.91</b>	<b>84.53</b>	<b>33.35</b>	39.91	<b>21.65</b>	<b>17.19</b>	<b>48.74</b>

Table 2. AP@0.15 scores for 10 out of 30 object categories [4] from the SUN RGB-D dataset.

Method	toilet	recycle bin	night stand	end table	drawer	computer	key board	table	chair	monitor	stool
PerspectiveNet[4]	81.22	37.68	35.16	19.77	1.28	1.24	<b>2.86</b>	44.12	40.42	1.14	<b>22.65</b>
ImVoxelNet	<b>84.53</b>	<b>52.20</b>	<b>46.29</b>	<b>25.31</b>	<b>6.05</b>	<b>2.71</b>	0.01	<b>51.17</b>	<b>65.94</b>	<b>19.82</b>	10.37

Method	lamp	dresser	picture	garbage bin	shelf	sofa chair	cabinet	sink	desk	book shelf	coffee table
PerspectiveNet[4]	13.14	<b>27.38</b>	0.00	22.42	0.97	51.86	1.70	<b>41.35</b>	20.19	8.29	28.80
ImVoxelNet	<b>17.19</b>	22.32	<b>0.82</b>	<b>33.35</b>	<b>4.00</b>	<b>54.61</b>	<b>7.90</b>	39.91	<b>31.91</b>	<b>21.65</b>	<b>36.48</b>

Method	box	sofa	white board	bed	pillow	paper	painting	cpu
PerspectiveNet[4]	1.64	62.35	0.02	<b>79.69</b>	11.36	0.00	0.17	<b>21.60</b>
ImVoxelNet	<b>3.29</b>	<b>63.89</b>	<b>0.95</b>	77.87	<b>14.65</b>	0.00	<b>0.53</b>	5.30

Table 3. AP@0.15 scores for 30 object categories [4] from the SUN RGB-D dataset.

Method	cabinet	bed	chair	sofa	table	door	window	book shelf	picture	counter	blinds
CooP[2]	10.47	57.71	15.21	36.67	31.16	0.14	0.00	3.81	0.00	27.67	<b>2.27</b>
T3DU[5]	11.39	59.03	15.98	43.95	35.28	0.36	0.16	5.26	<b>0.24</b>	<b>33.51</b>	0.00
ImVoxelNet	<b>19.24</b>	<b>79.17</b>	<b>63.07</b>	<b>60.59</b>	<b>51.14</b>	<b>0.74</b>	<b>0.18</b>	<b>16.37</b>	0.14	14.89	0.26

Method	desk	shelves	curtain	dresser	pillow	mirror	floor mat	clothes	books	fridge	tv
CooP[2]	19.90	2.96	1.35	15.98	2.53	<b>0.47</b>	-	0.00	<b>3.19</b>	21.50	5.20
T3DU[5]	23.65	4.96	2.68	19.20	2.99	0.19	-	0.00	1.30	20.68	4.44
ImVoxelNet	<b>31.20</b>	<b>5.47</b>	<b>3.34</b>	<b>35.45</b>	<b>11.01</b>	0.22	-	<b>1.40</b>	0.13	<b>23.28</b>	<b>12.41</b>

Method	paper	towel	shower curtain	box	white board	person	night stand	toilet	sink	lamp	bathtub
CooP[2]	0.20	2.14	<b>20.00</b>	2.59	0.16	20.96	11.36	42.53	15.95	3.28	24.71
T3DU[5]	<b>0.41</b>	<b>2.20</b>	<b>20.00</b>	2.25	0.43	23.36	6.87	48.37	14.40	3.46	27.85
ImVoxelNet	0.00	1.92	0.00	<b>2.71</b>	<b>1.17</b>	<b>42.02</b>	<b>38.38</b>	<b>77.28</b>	<b>45.12</b>	<b>13.27</b>	<b>43.59</b>

Method	bag	wall	floor	ceiling
CooP[2]	1.53	-	-	-
T3DU[5]	<b>2.27</b>	-	-	-
ImVoxelNet	0.53	-	-	-

Table 4. AP@0.15 scores for 37 object categories [5] from the SUN RGB-D dataset.

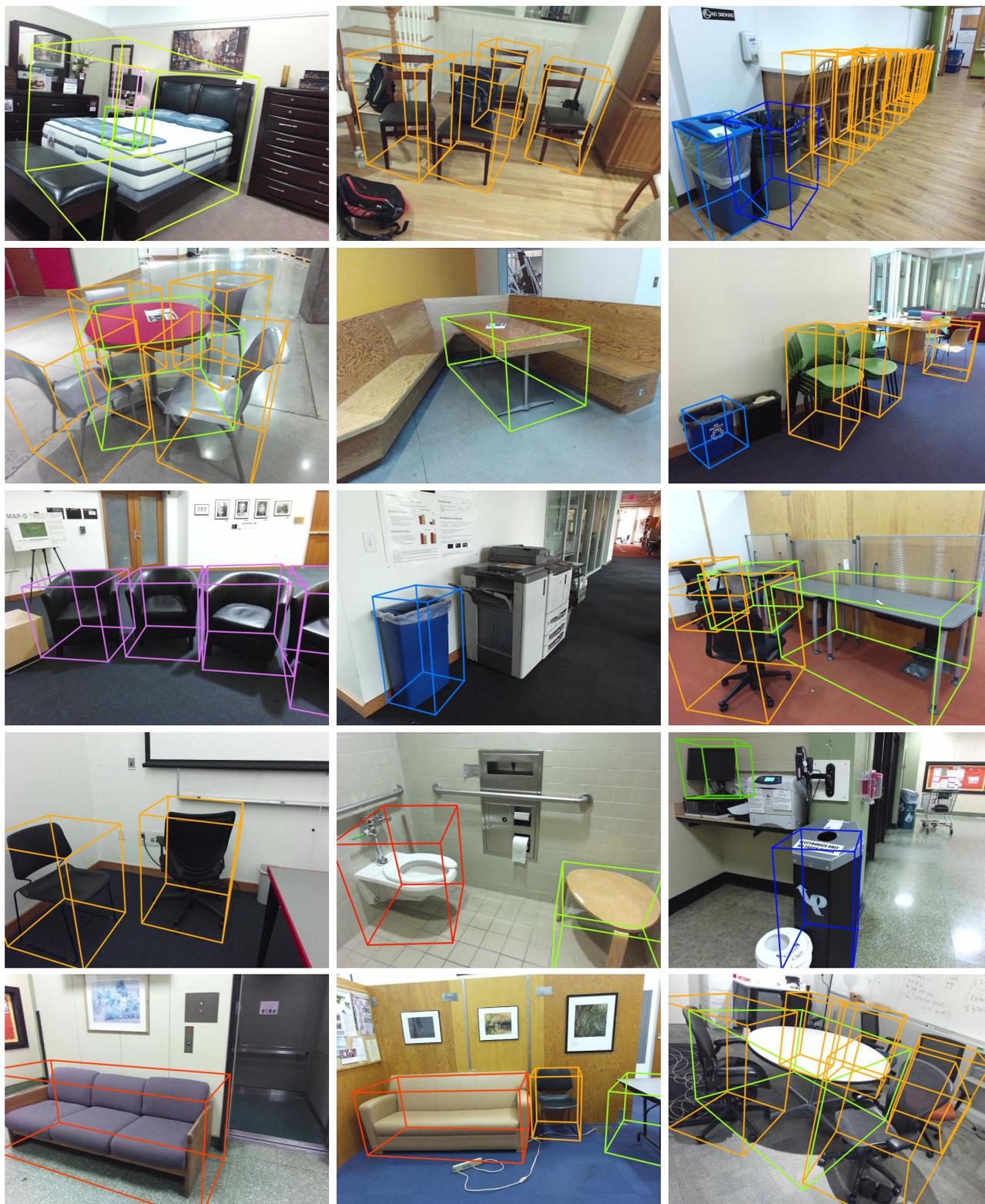
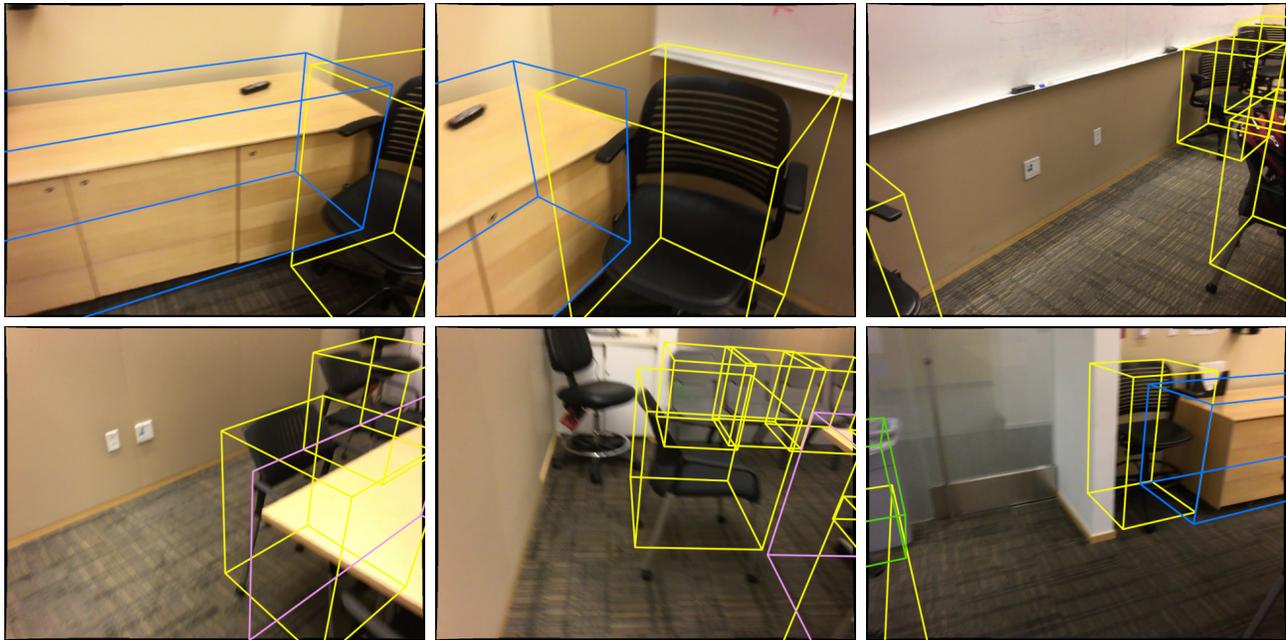
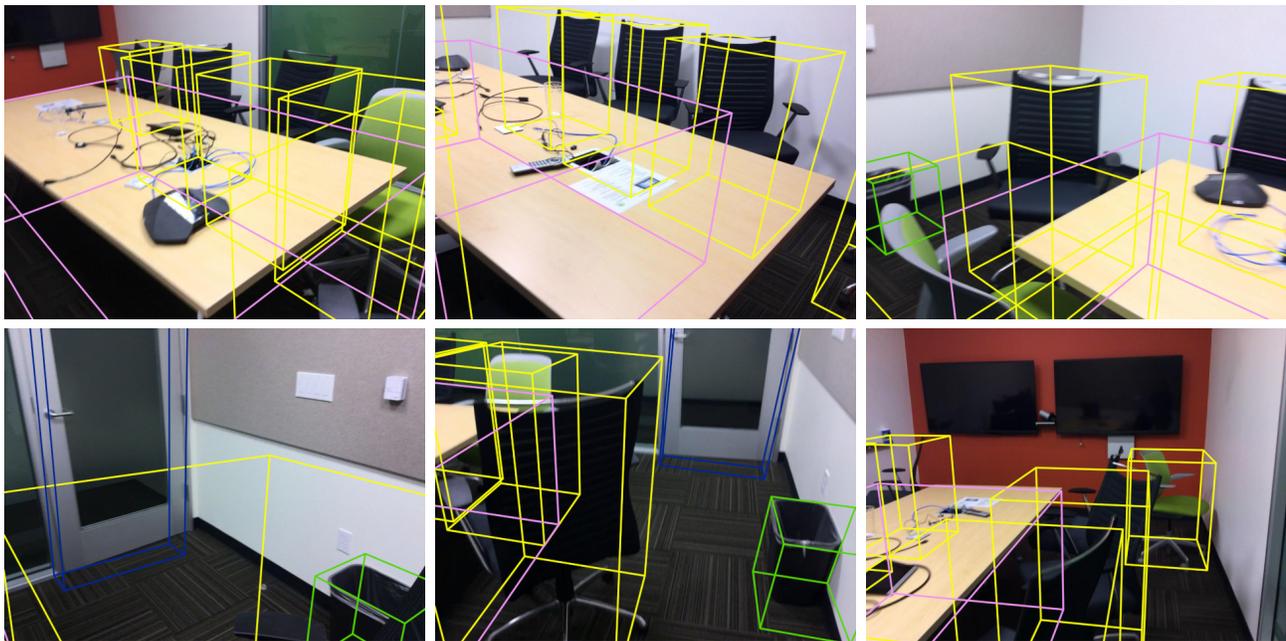


Figure 1. Objects detected on the monocular images from the validation subset of the SUN RGB-D dataset.



a) Scene 0169\_00.

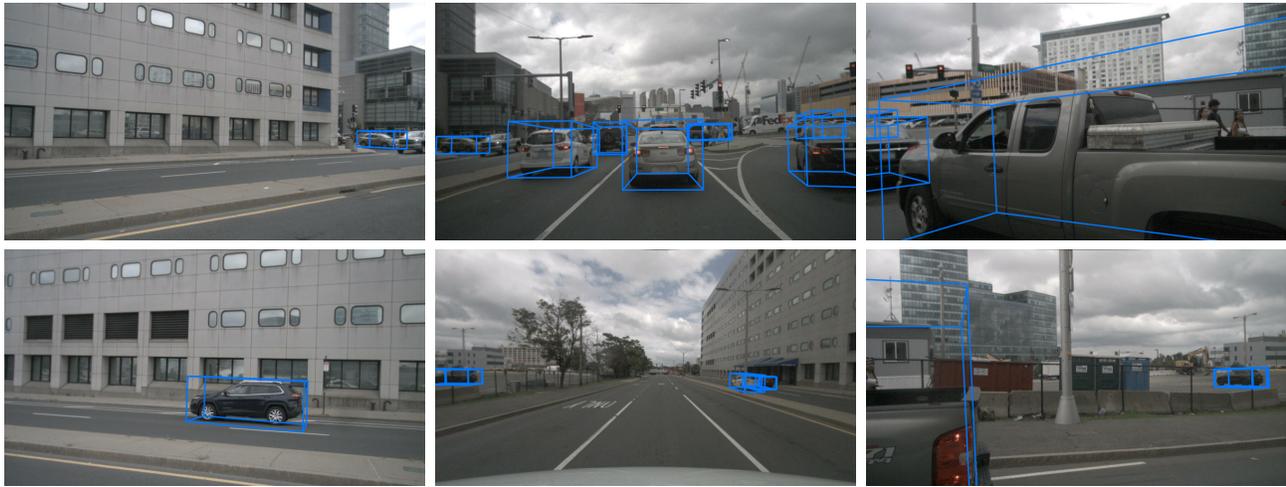


b) Scene 0575\_00.

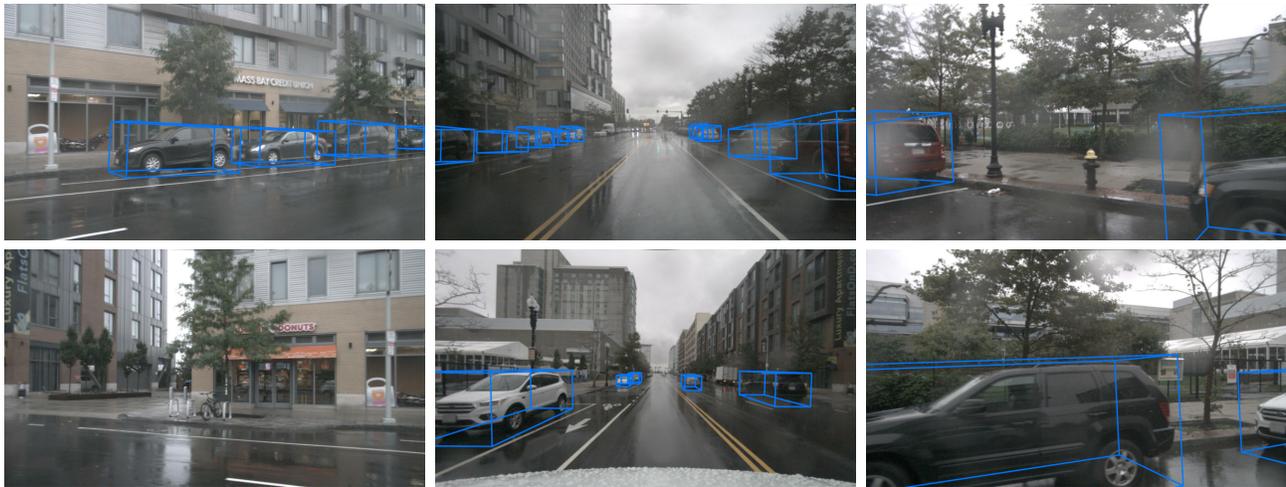
Figure 2. Objects detected on the multi-view inputs from the validation subset of the ScanNet dataset.



Figure 3. Cars detected on the monocular images from the validation subset of the KITTI dataset.



a) Scene *n008-2018-08-01-15-16-36-0400\_15331512526*.



b) Scene *n008-2018-09-18-15-12-01-0400\_15372981046*.

Figure 4. Cars detected in the images of two scenes from the validation subset of the nuScenes dataset. The predictions were obtained in multi-view settings. The first two rows correspond to the first scene, and the last two rows correspond to another one. For each scene, the upper row consists of images taken with a front-left, front, and front-right camera (from left to right). The second row contains images taken with a back-left, back, and back-right camera, respectively.

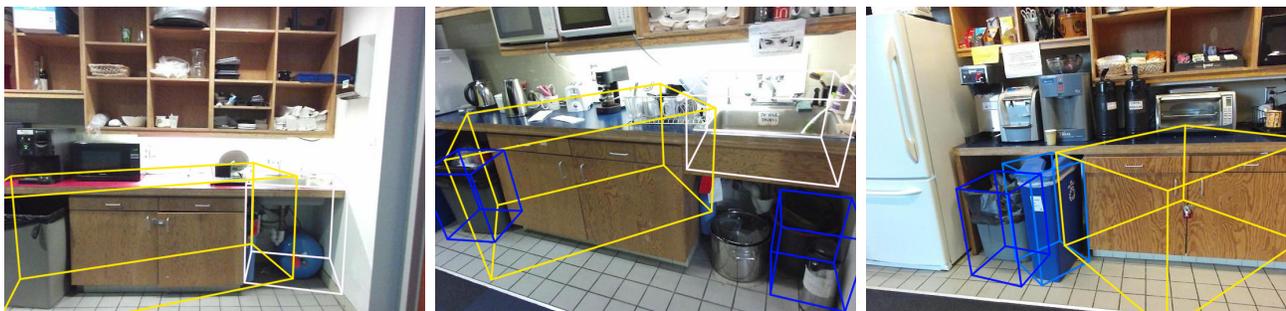


Figure 5. Examples of the detection failures for images from the validation subset of the SUN RGB-D dataset. These examples depict typical error cases: small objects of *sink*, *garbage bin* and *recycle bin* categories are detected quite precisely, but rotation angles for large object such as *cabinet* are estimated poorly.

## References

- [1] W. Choi, Y.-W. Chao, C. Pantofaru, and S. Savarese. Understanding indoor scenes using 3d geometric phrases. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 33–40, 2013.
- [2] S. Huang, S. Qi, Y. Xiao, Y. Zhu, Y. N. Wu, and S.-C. Zhu. Cooperative holistic scene understanding: Unifying 3d object, layout, and camera pose estimation. *arXiv preprint arXiv:1810.13049*, 2018.
- [3] S. Huang, S. Qi, Y. Zhu, Y. Xiao, Y. Xu, and S.-C. Zhu. Holistic 3d scene parsing and reconstruction from a single rgb image. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 187–203, 2018.
- [4] S. Huang, Y. Chen, T. Yuan, S. Qi, Y. Zhu, and S.-C. Zhu. Perspectivenet: 3d object detection from a single rgb image via perspective points. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [5] Y. Nie, X. Han, S. Guo, Y. Zheng, J. Chang, and J. J. Zhang. Total3dunderstanding: Joint layout, object pose and mesh reconstruction for indoor scenes from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 55–64, 2020.
- [6] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 918–927, 2018.
- [7] C. R. Qi, O. Litany, K. He, and L. J. Guibas. Deep hough voting for 3d object detection in point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9277–9286, 2019.
- [8] C. R. Qi, X. Chen, O. Litany, and L. J. Guibas. Imvotenet: Boosting 3d object detection in point clouds with image votes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4404–4413, 2020.
- [9] Z. Zhang, B. Sun, H. Yang, and Q. Huang. H3dnet: 3d object detection using hybrid geometric primitives. In *European Conference on Computer Vision*, pages 311–329. Springer, 2020.