



# From Points to Multi-Object 3D Reconstruction

Francis Engelmann<sup>1</sup> Konstantinos Rematas<sup>2</sup> Bastian Leibe<sup>1</sup> Vittorio Ferrari<sup>2</sup>

<sup>1</sup>RWTH Aachen University, Germany    <sup>2</sup>Google Research, Zurich

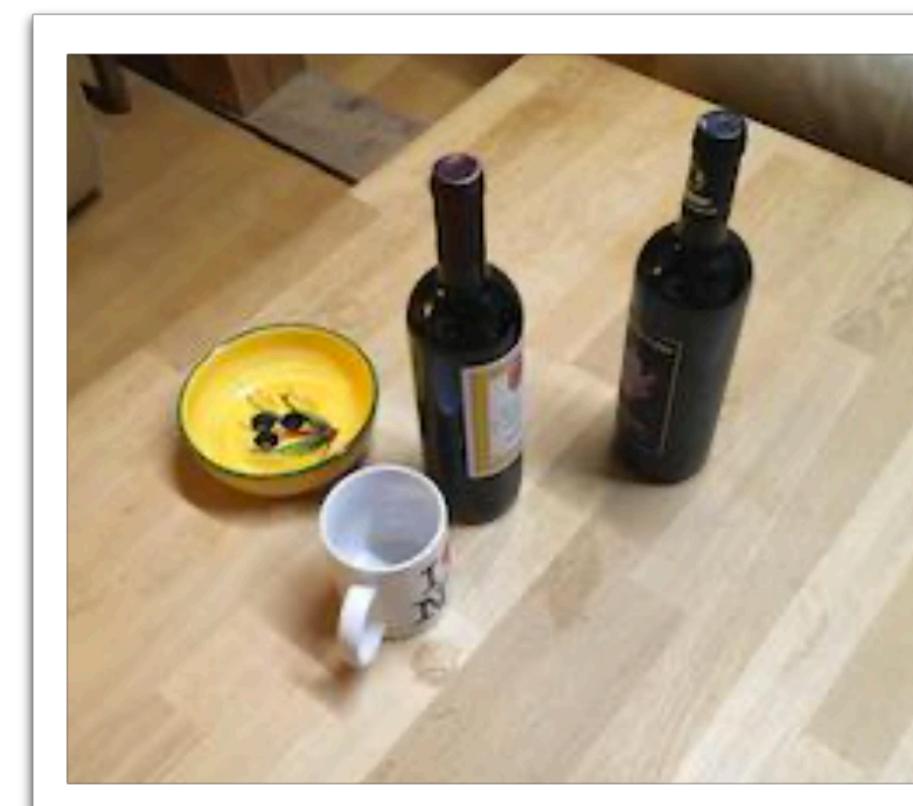
<https://francisengelmann.github.io/points2objects/>



## Abstract

We propose a method to detect and reconstruct multiple 3D objects from a single RGB image. The key idea is to optimize for detection, alignment and shape jointly over all objects in the RGB image, while focusing on realistic and physically plausible reconstructions. To this end, we propose a key-point detector that localizes objects as center points and directly predicts all multi-object properties, including 9-DoF bounding boxes and 3D shapes — all in a single forward pass.

## Task



Points2Objects

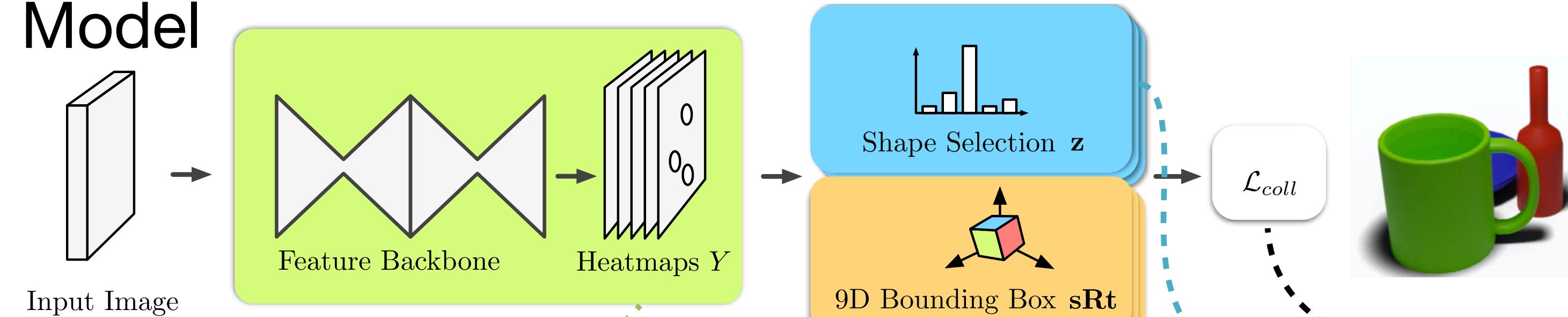


Input: Single Image

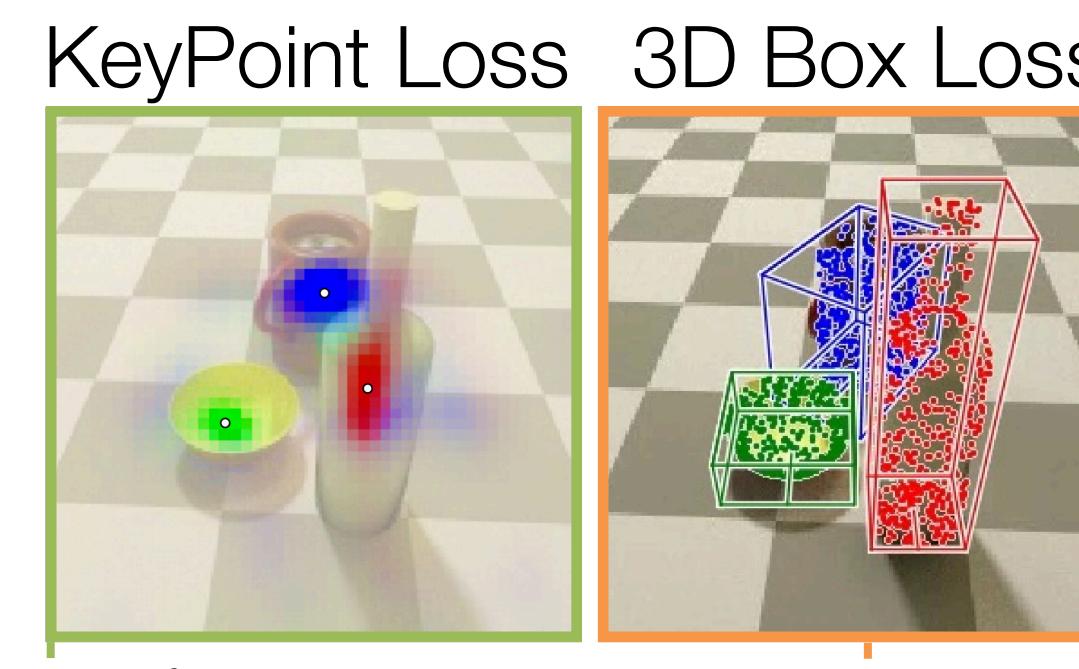
## Contributions

- Fully holistic multi-object 3D scene reconstruction based on CenterNet [1] in a single-stage network from a single input RGB image.
- Our reconstruction is formulated as a shape-selection problem (1-of-K classification) implemented using our novel “soft target labels” relying on geometric similarities between exemplar 3D shapes.
- Our collision loss encourages non-intersecting reconstructions and CAD representations guarantee physically plausible and realistic shapes.
- We present a 9-DoF pose estimation study showing that jointly optimizing rotation and translation improves over individual optimization in our setup.

## Model

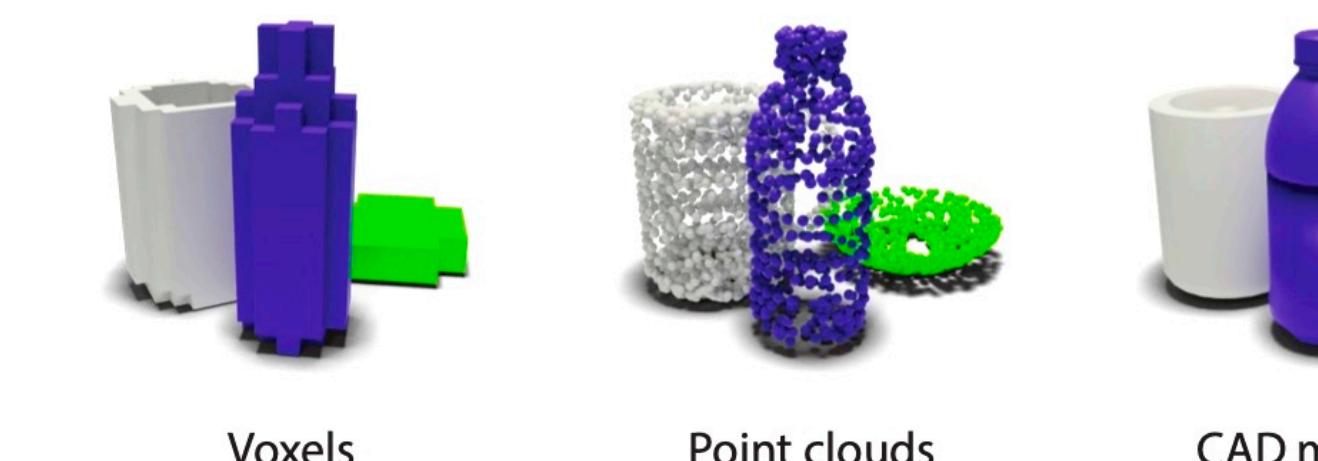


## Losses



$\mathcal{L}_{key} = \text{as in CenterNet [1]}$   
 $\mathcal{L}_{Rt} = \sum_{i=1}^M \sum_{x \in P^i} \|[\mathbf{R}|t]^i \mathbf{x} - [\hat{\mathbf{R}}|\hat{t}]^i \mathbf{x}\|_2^2$

## Representation-agnostic Shapes



## Evaluation

### Estimating 9-DoF Poses - Study

9-DoF Bounding Box	3D mAP:	@ 0.5	@ 0.25
$\mathcal{L}_{binR} + \mathcal{L}_{offR} + \mathcal{L}_t$ (as in [1])	43.3	75.0	
$\mathcal{L}_M + \mathcal{L}_t$ (directly regress rotation matrix M)	44.8	77.0	
$\mathcal{L}_R + \mathcal{L}_t$ (add SVD for orthogonal rotation R [5])	46.8	77.2	
$\mathcal{L}_{Rt}$ (ours)	<b>48.6</b>	<b>77.2</b>	

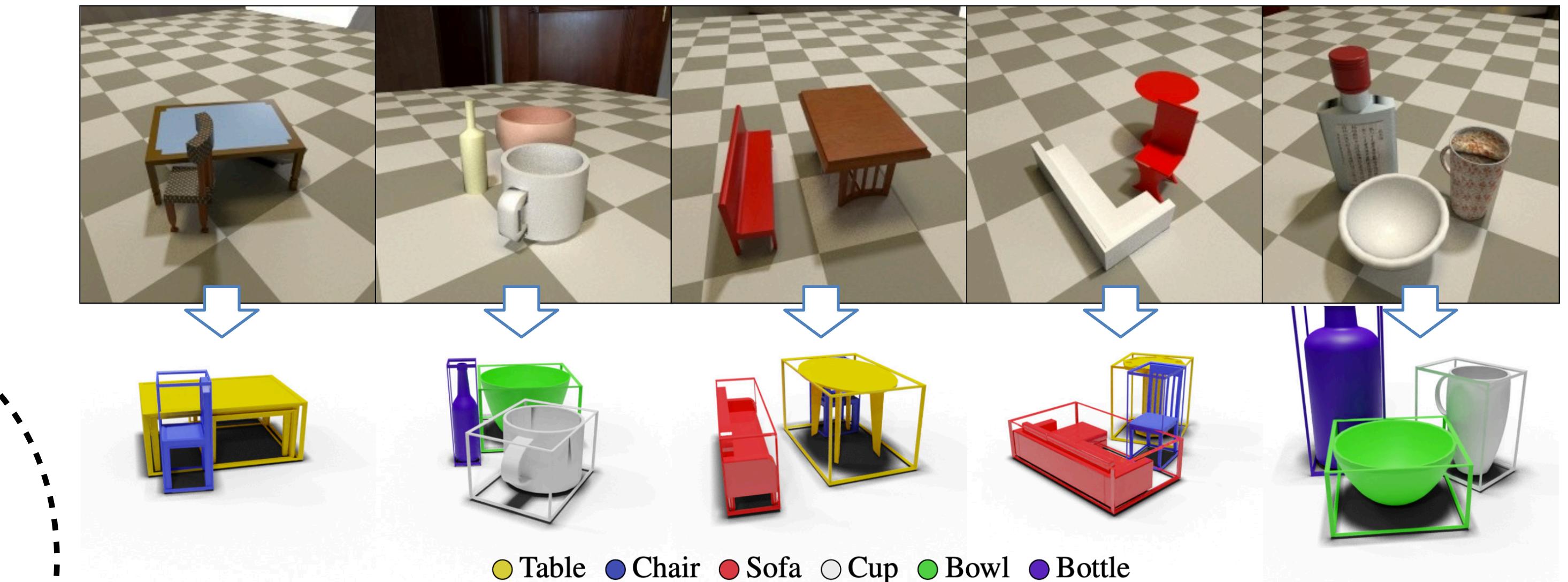
### Effect of collision loss

	Num. Collisions
Without collision loss	4116
With collision loss	<b>1627</b> $\downarrow -60.5\%$

### Shape Estimation: Hard vs. Soft Labels

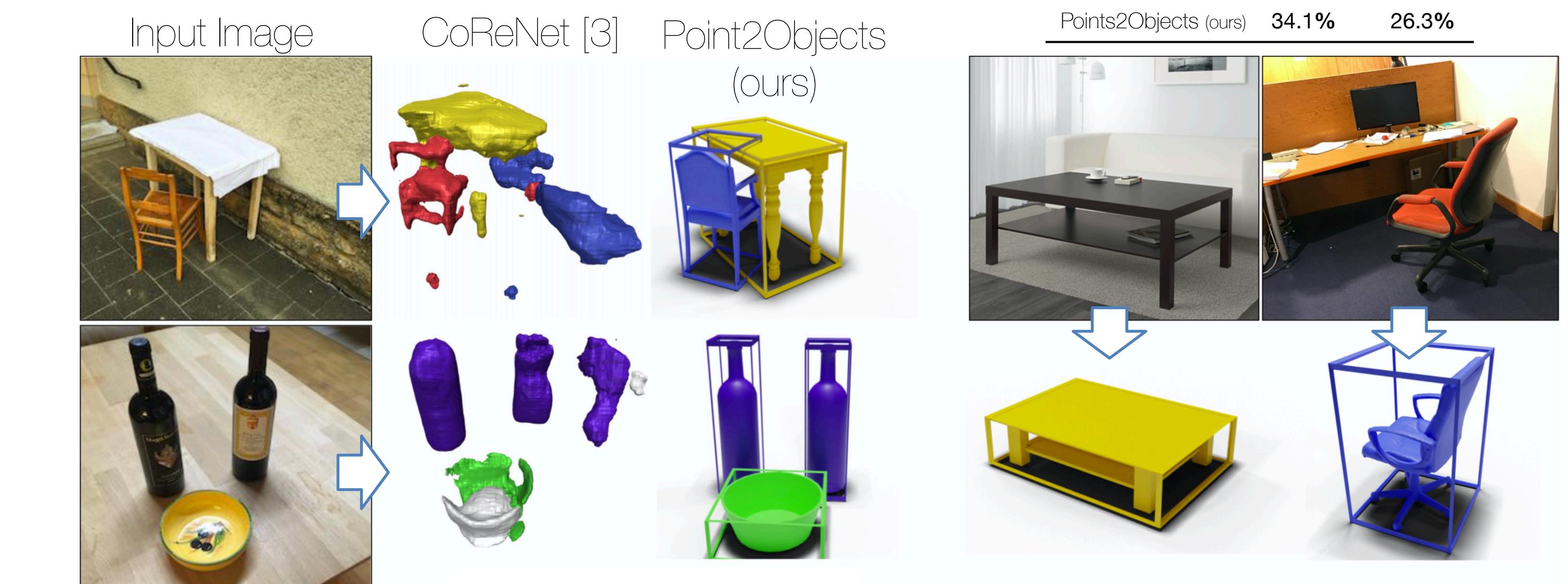
Shape Estimation	Abs. 3D IoU:	mean	global
$\mathcal{L}'_z$ Hard-Labels (as in [2])	32.2	40.3	
$\mathcal{L}_z$ Soft-Labels (ours)	<b>36.4</b>	<b>44.7</b>	

## Results on synthetic images



## Results on real images

Casual photos from mobile phone  
Generalization from synthetic to real data



## References

- Xingyi Zhou, Dequan Wang, Philipp Krähenbühl. “Objects as Points” ArXiv 2019.
- Maxim Tatarchenko, Stephan R. Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, Thomas Brox “What Do Single-view 3D Reconstruction Networks Learn?” In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- Stefan Popov, Pablo Bausat, Vittorio Ferrari “CoReNet: Coherent 3D Scene Reconstruction from a Single RGB Image” In IEEE European Conference on Computer Vision (ECCV), 2020.
- Xingyuan Sun\*, Jiajun Wu\*, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue, Joshua B. Tenenbaum, and William T. Freeman “Pix3D: Dataset and Methods for Single-Image 3D Shape Modeling” In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.
- J. Levinson, Carlos Esteves, Kefan Chen, Noah Snavely, A. Kanazawa, Afshin Rostamizadeh, and A. Makadia. “An Analysis of SVD for Deep Rotation Estimation”. In Neural Information Processing Systems (NeurIPS), 2020.

Pix3D [4] Single object dataset	Splits	S <sub>1</sub>	S <sub>2</sub>
CoReNet [3]		33.3%	23.6%
Point2Objects (ours)		<b>34.1%</b>	<b>26.3%</b>

