

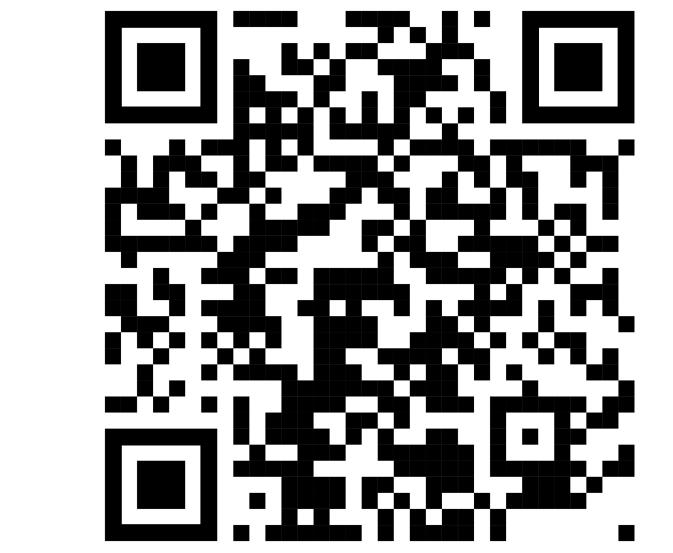


From Points to Multi-Object 3D Reconstruction

Francis Engelmann¹ Konstantinos Rematas² Bastian Leibe¹ Vittorio Ferrari²

¹RWTH Aachen University, Germany ²Google Research, Zurich

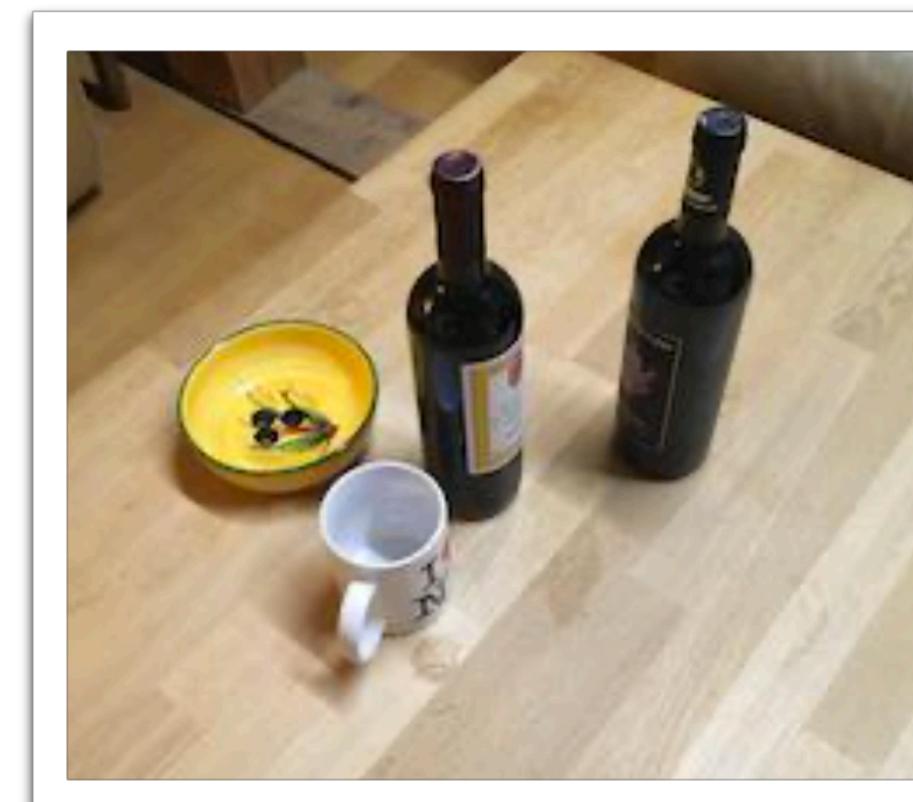
<https://francisengelmann.github.io/points2objects/>



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.

The Task



Input: Single Image

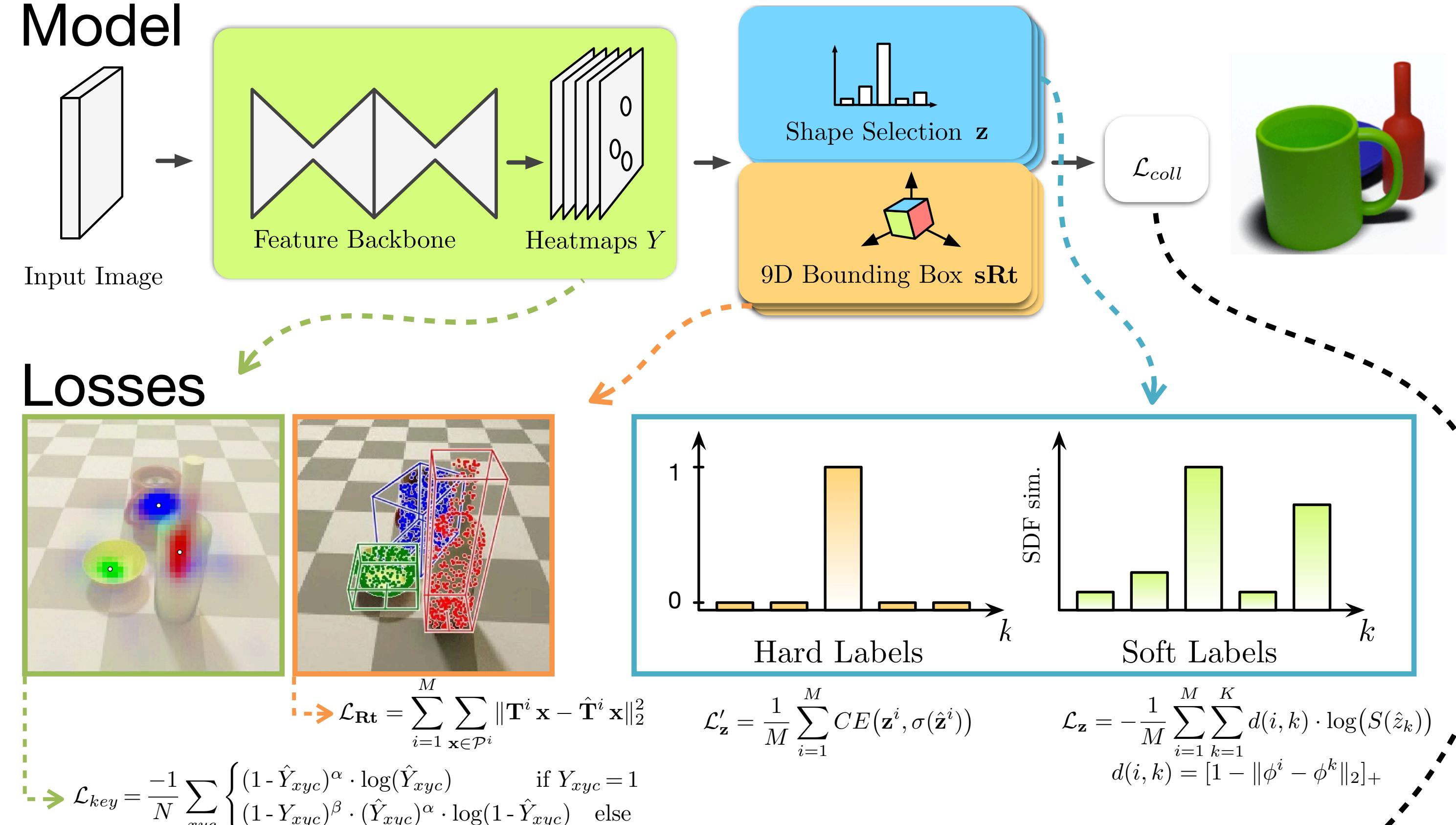


Output: 3D Reconstructions of Multiple Objects

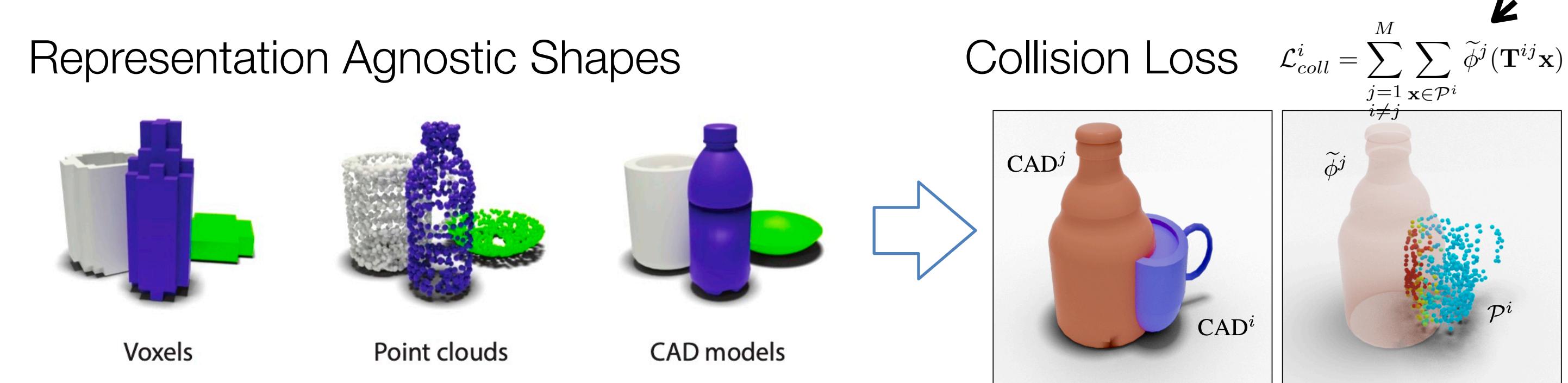
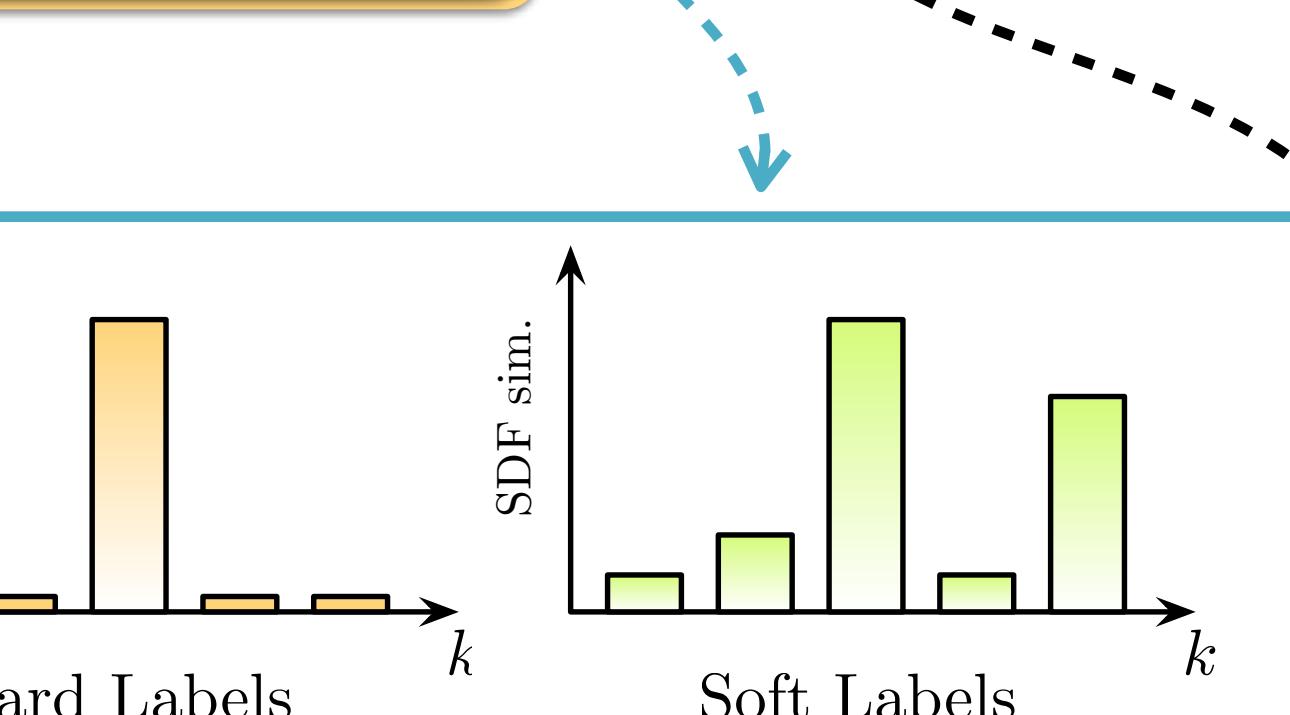
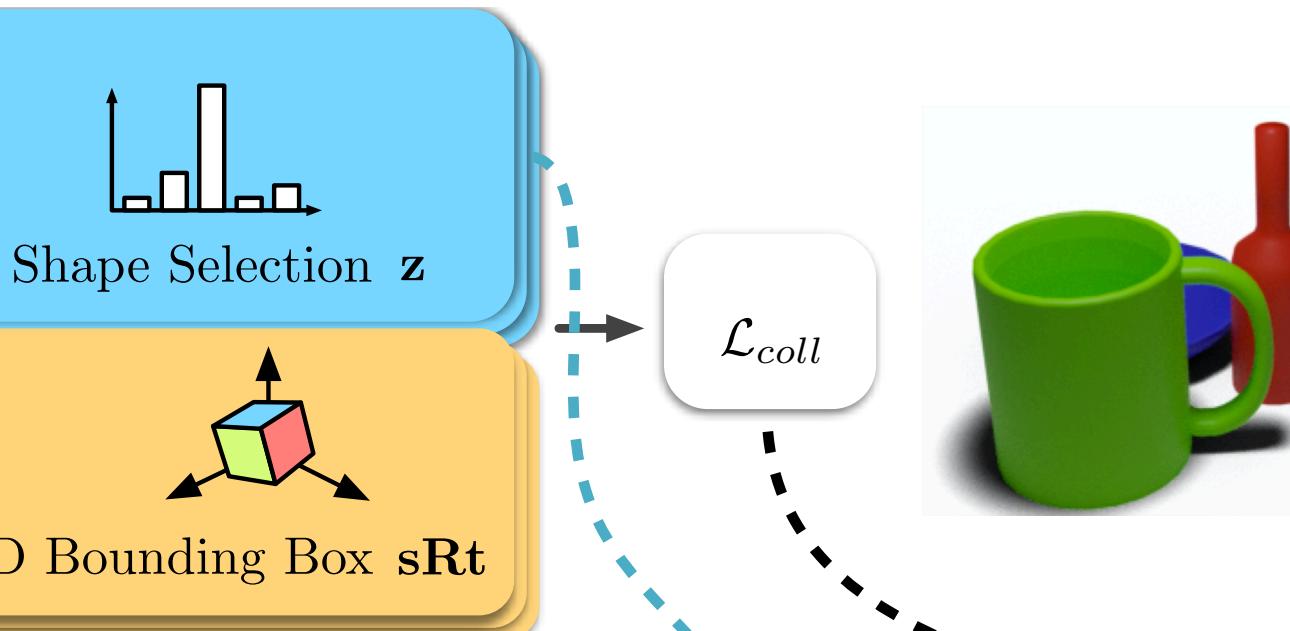
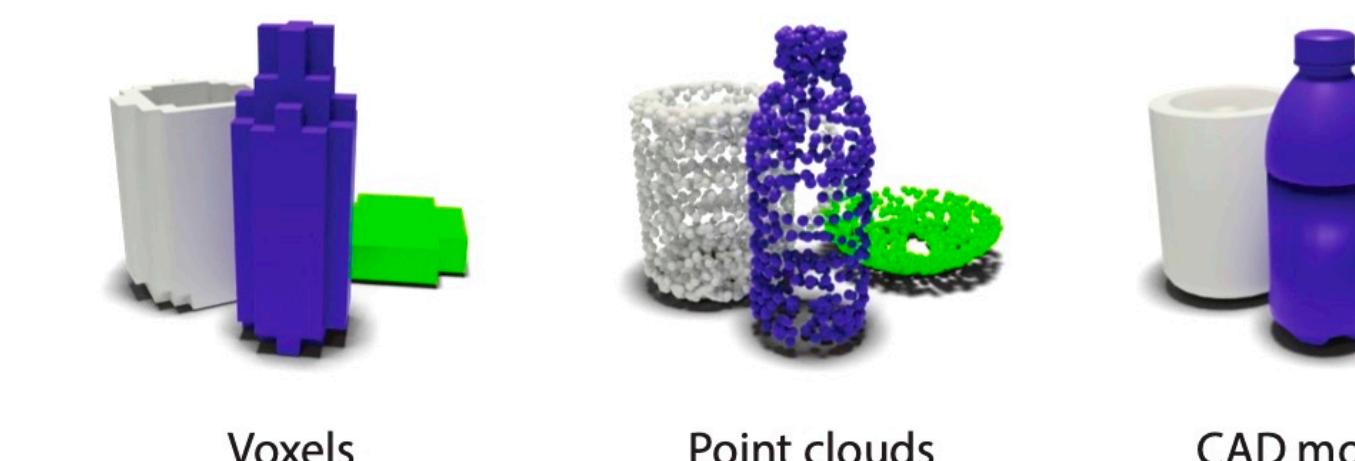
The Contributions

- 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-based representations guarantee physically plausible and realistic shapes.

Model



Representation Agnostic Shapes



Evaluation

Optimizing Bounding Boxes - Study

9-DoF Bounding Box	3D mAP:	@ 0.5	@ 0.25
$\mathcal{L}_{\text{binR}} + \mathcal{L}_{\text{offR}} + \mathcal{L}_t$ (as in [1])	43.3	75.0	
$\mathcal{L}_M + \mathcal{L}_t$	44.8	77.0	
$\mathcal{L}_R + \mathcal{L}_t$ (with SVD \rightarrow valid rotation)	46.8	77.2	
\mathcal{L}_{Rt} (ours)	48.6	77.2	

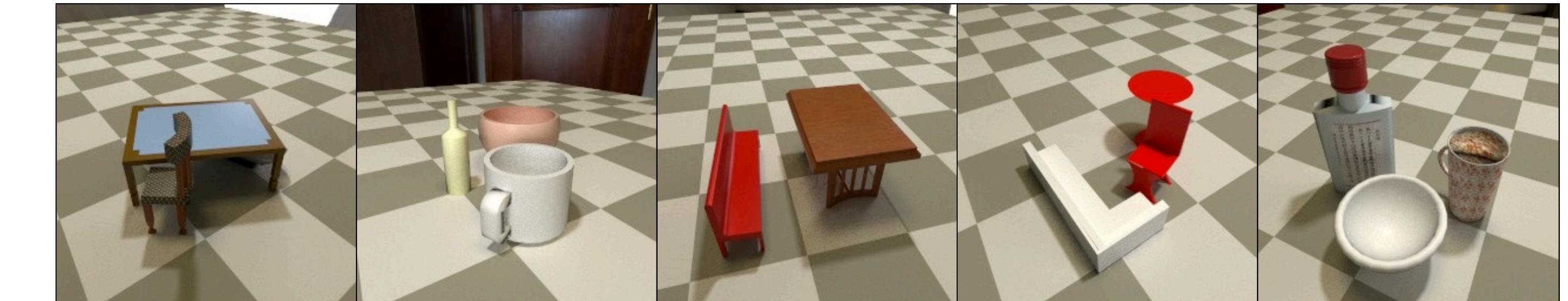
Collision Loss Effect

mIV	Num. Collisions
1168.8	4116
794	1627 $\downarrow 60.5\%$

Shape Estimation: Hard vs. Soft Labels

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

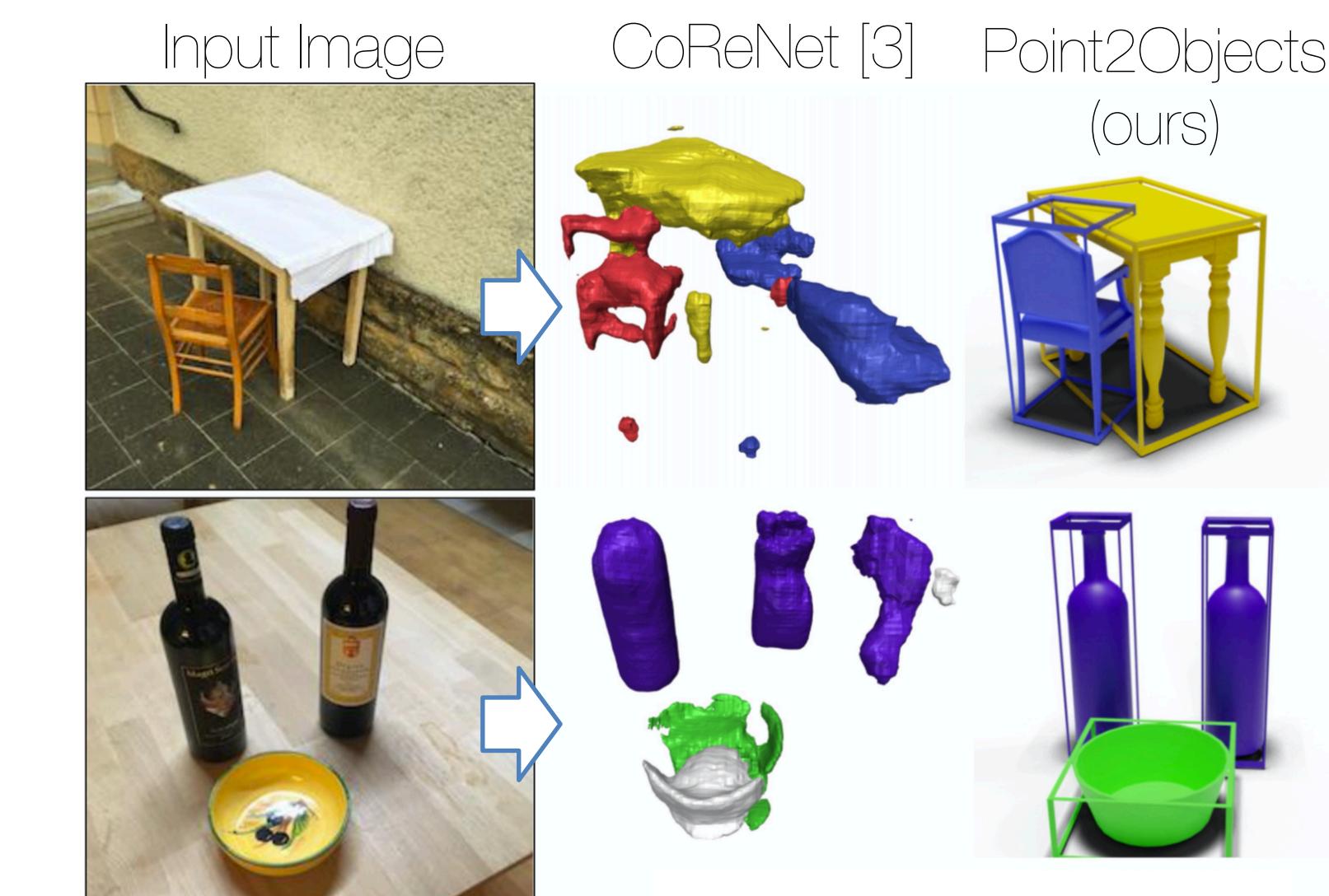
Results on synthetic images



Yellow Table, Blue Chair, Red Sofa, White Cup, Green Bowl, Purple Bottle

Results on real images

Casual photos from mobile phone
Generalization from synthetic to real data



Pix3D [4]
Single object dataset



Pix3D [4]
Single object dataset

References

- [1] Xingyi Zhou, Dequan Wang, Philipp Krähenbühl. “Objects as Points” ArXiv 2019.
- [2] 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.
- [3] Stefan Popov, Pablo Bauszat, Vittorio Ferrari “CoReNet: Coherent 3D Scene Reconstruction from a Single RGB Image” In IEEE European Conference on Computer Vision (ECCV), 2020.
- [4] 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), 2017.