

# Learning Semantic Perception for Cluttered Bin Picking

**Sven Behnke and Max Schwarz**

University of Bonn

Computer Science Institute VI

Autonomous Intelligent Systems



# Bin Picking

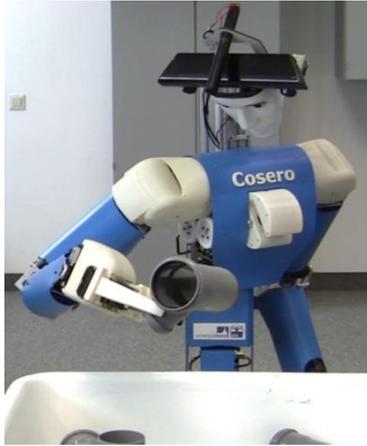
- Removing items from containers and shelves
- Still often performed by humans
- Difficulties include
  - Item variability
  - Problematic material properties
  - Articulation of objects
  - Lacking grasp affordances
  - Chaotic storage
  - Inaccessibility



[Amazon]

# Our Past Experience

ActReMa



STAMINA



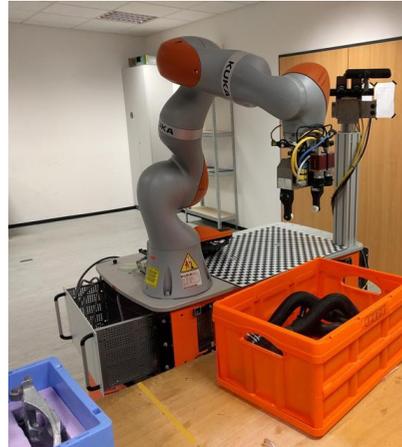
Amazon Picking



EuRoC C1



EuRoC C2



# Amazon Robotics Challenge 2017

- Quickly learn novel objects
- Design own storage system

Sensor setup

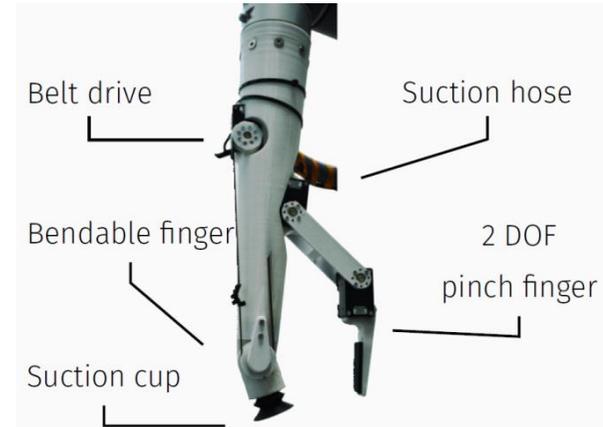
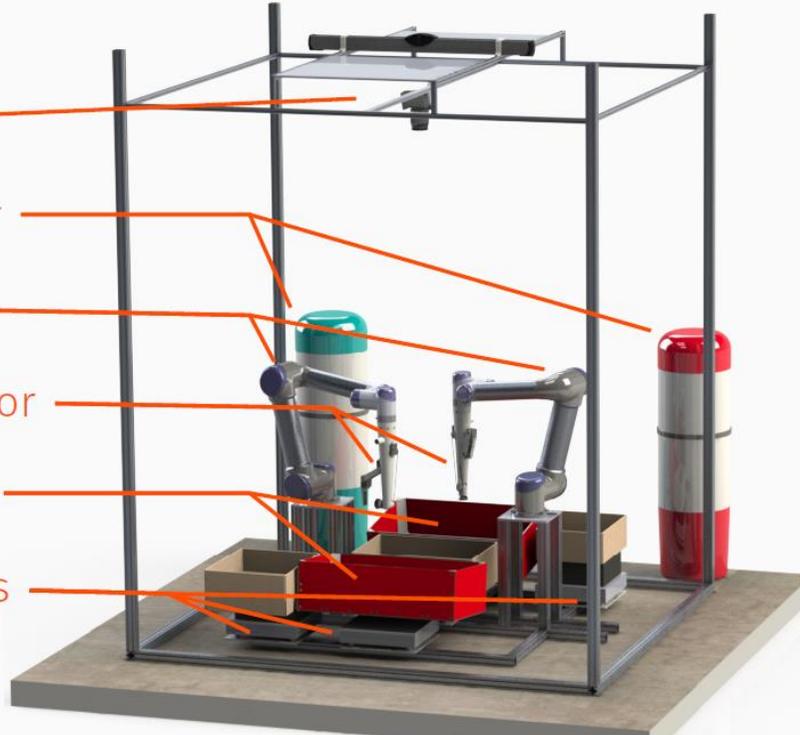
Vacuum cleaner

6 DOF UR5 arm

3 DOF endeffector

Storage system

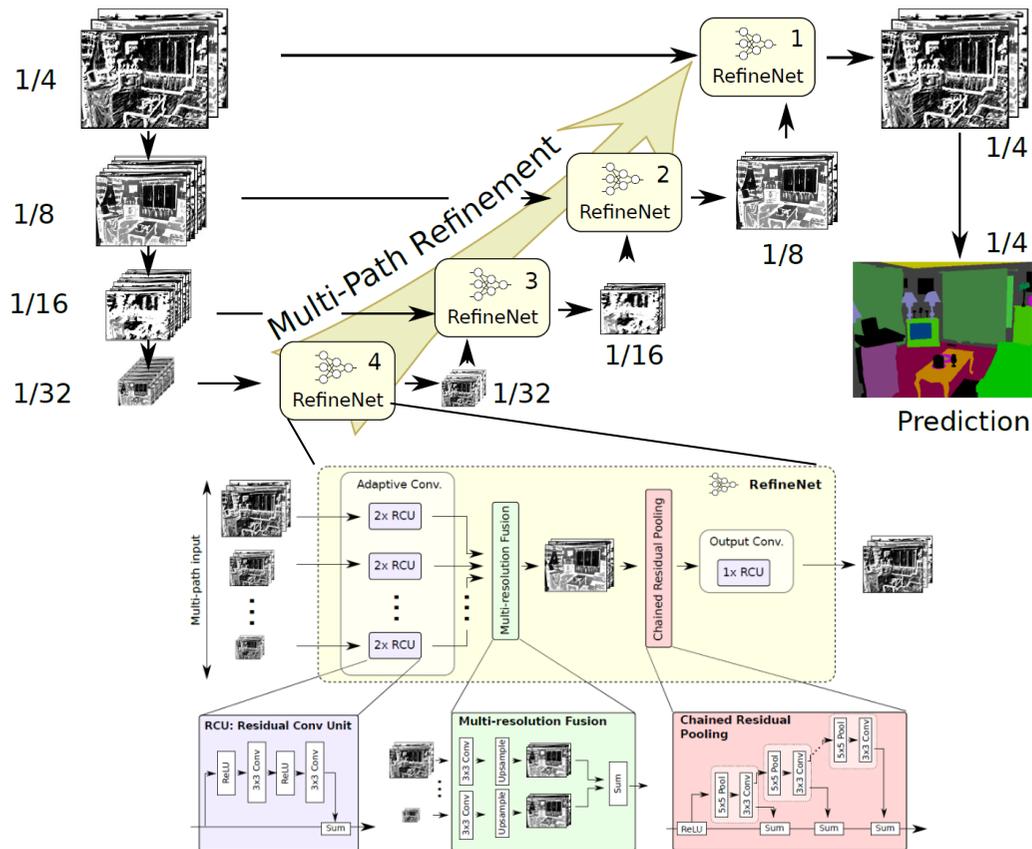
Industrial scales



[Schwarz et al. ICRA 2018]

# RefineNet for Semantic Segmentation

- Scene represented as feature hierarchy
- Coarse-to-fine semantic segmentation
- Combine higher-level features with missing details



[Lin et al. CVPR 2017]

# The Data Problem

- Deep Learning in robotics (still) suffers from shortage of available examples
- We address this problem in two ways:

## 1. Generating data:

Automatic data capture,  
online mesh databases,  
scene synthesis

## 2. Improving generalization:

Object-centered models,  
deformable registration,  
transfer learning,  
semi-supervised learning



# Object Capture and Scene Rendering

## ■ Turntable + DLSR camera

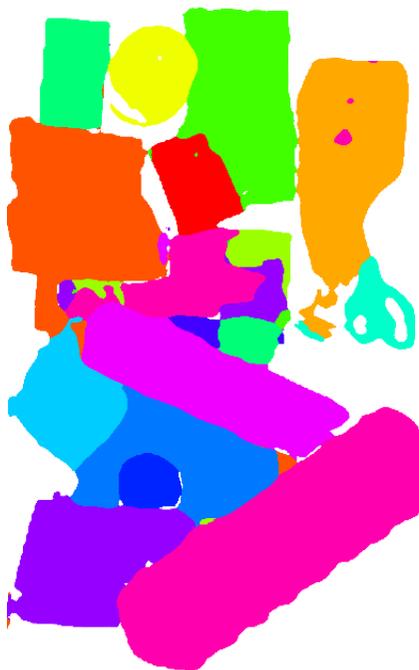


## ■ Rendered scenes



[Schwarz et al. ICRA 2018]

# ARC 2017 Perception Example



- bronze\_wire\_cup  
conf: 0.749401
- irish\_spring\_soap  
conf: 0.811500
- playing\_cards  
conf: 0.813761
- w\_aquarium\_gravel  
conf: 0.891001
- crayons  
conf: 0.422604
- reynolds\_wrap  
conf: 0.836467
- paper\_towels  
conf: 0.903645
- white\_facecloth  
conf: 0.895212
- hand\_weight  
conf: 0.928119
- robots\_everywhere  
conf: 0.930464



- mouse\_traps  
conf: 0.921731
- windex  
conf: 0.861246
- q-tips\_500  
conf: 0.475015
- fiskars\_scissors  
conf: 0.831069
- ice\_cube\_tray  
conf: 0.976856

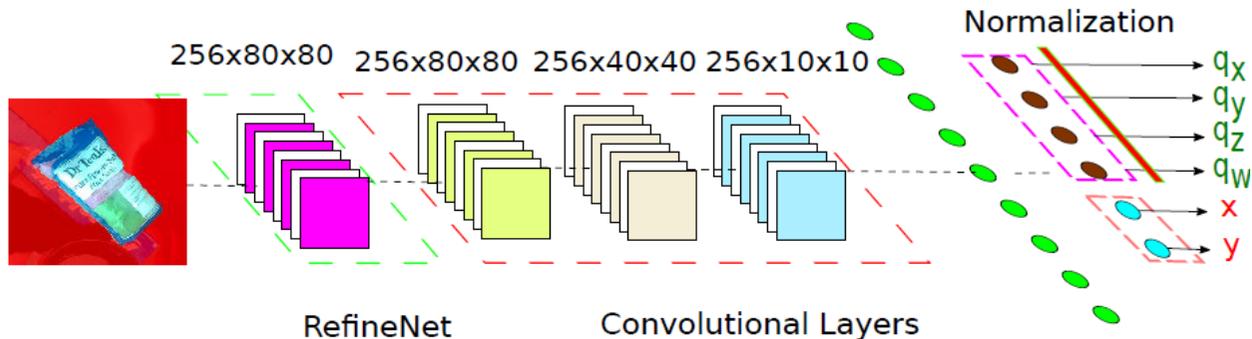
# Amazon Robotics Challenge 2017 Final



[Schwarz et al. ICRA 2018]

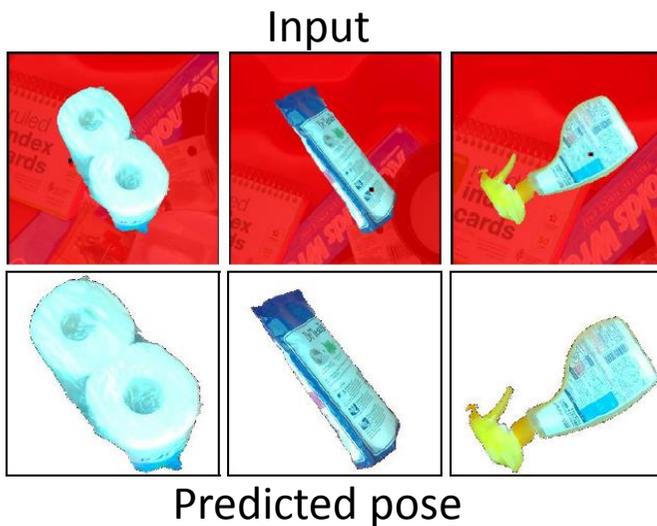
# Object Pose Estimation

- Cut out individual segments
- Use upper layer of RefineNet as input
- Predict pose coordinates



Object	Translation [pixel] <sup>1</sup>		Orientation [°]	
	train	val	train	val
Browns brush	10.3	11.4	7.7	10.3
Epsom salts	11.2	12.5	7.4	10.5
Hand weight	9.6	10.4	2.1	2.6
Reynolds wrap	11.6	11.8	6.3	9.8
Utility brush	12.5	13.6	6.9	10.9

<sup>1</sup> Relative to the 320×320 crop centered on the object.



# Transfer of Manipulation Skills

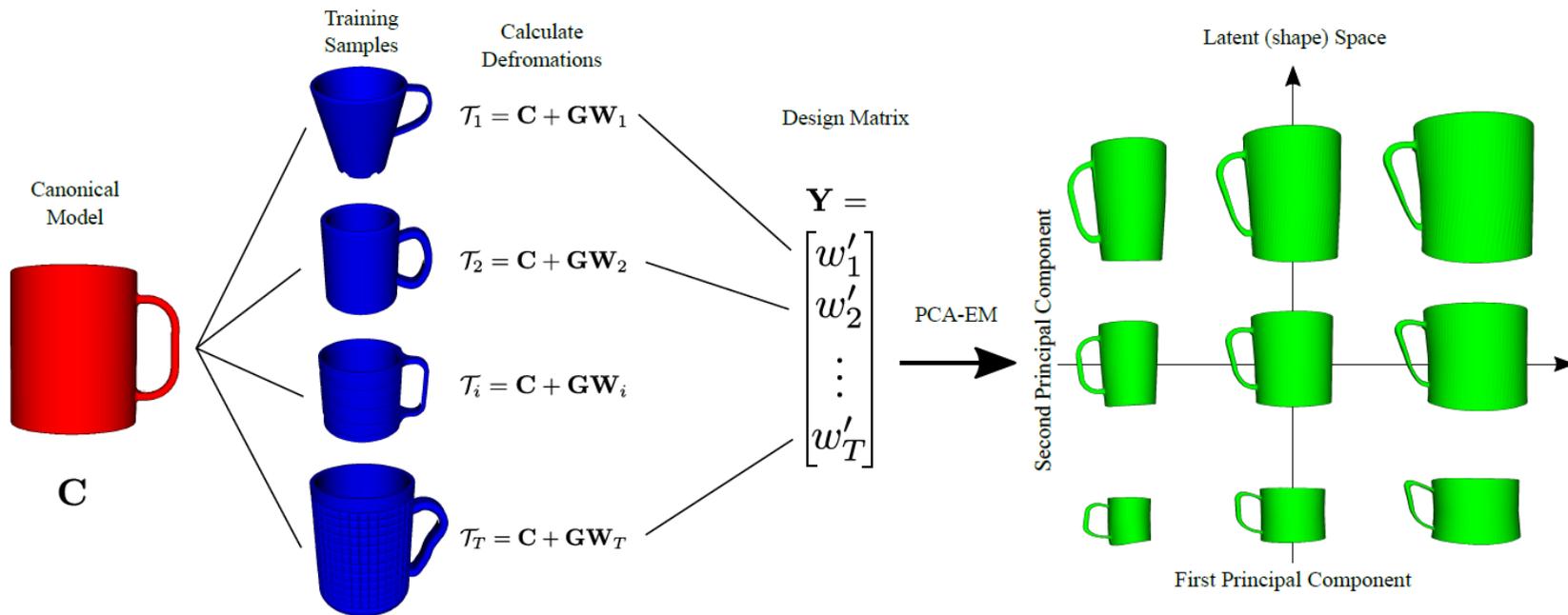


Knowledge  
Transfer

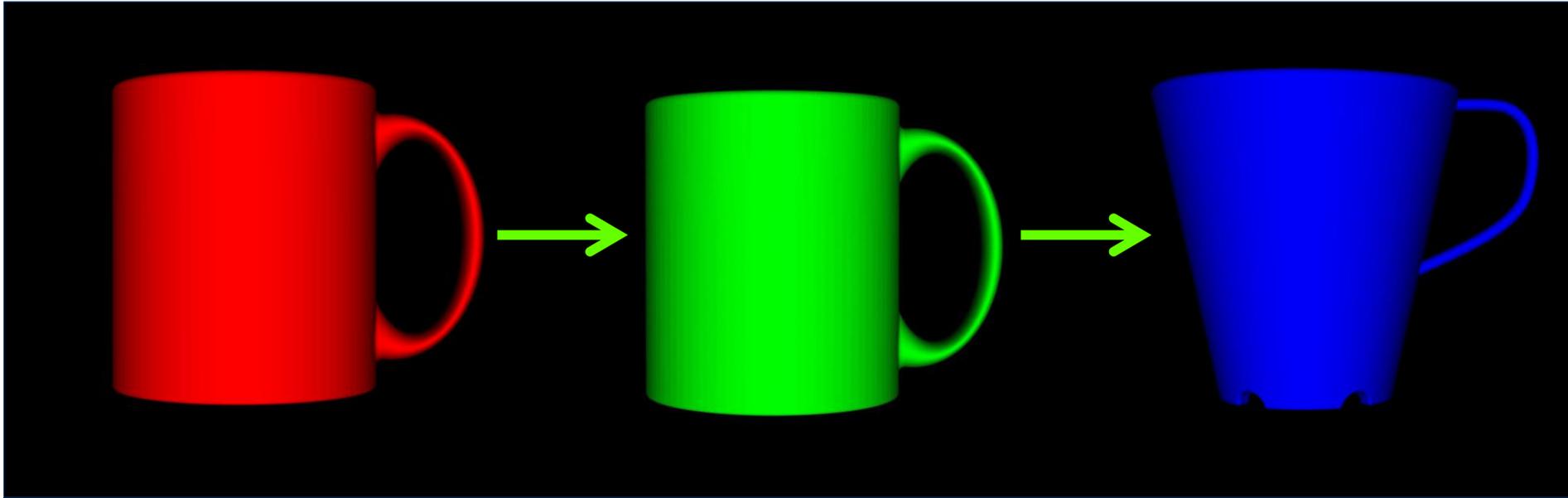


# Learning a Latent Shape Space

- Non-rigid registration of instances to canonical model
- Principal component analysis of deformations

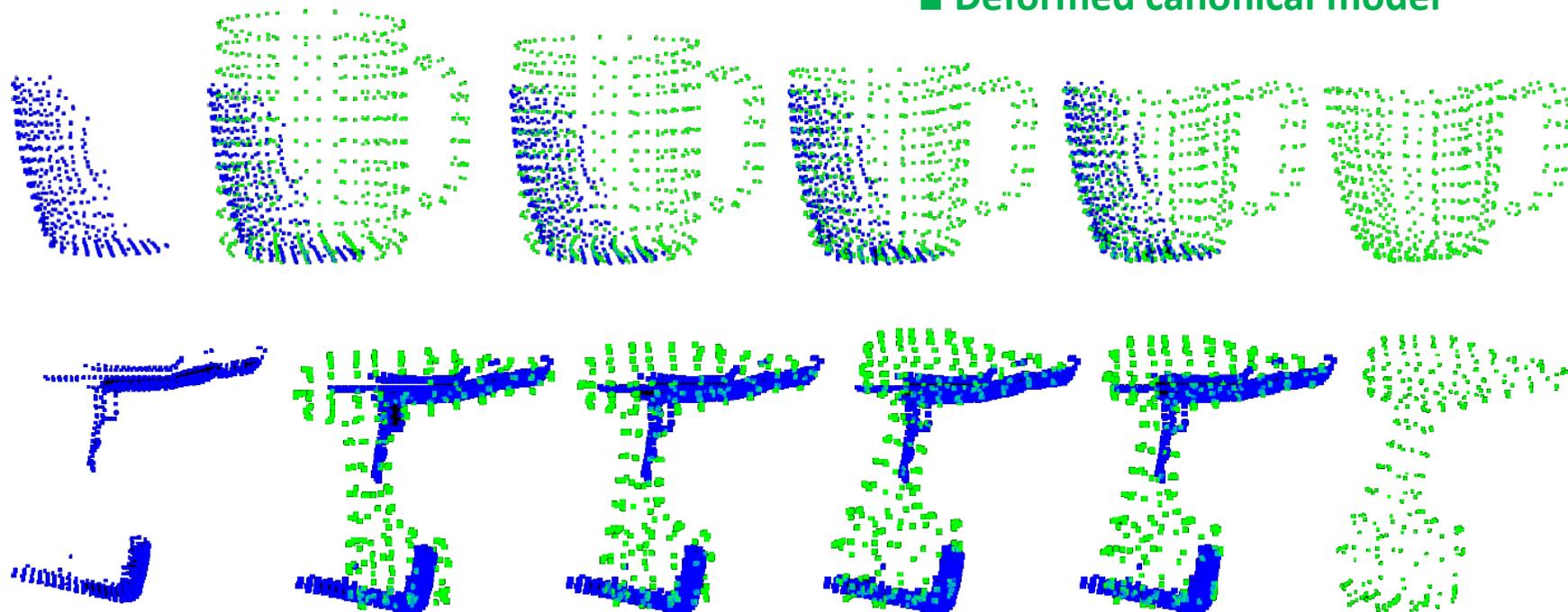


# Interpolation in Shape Space



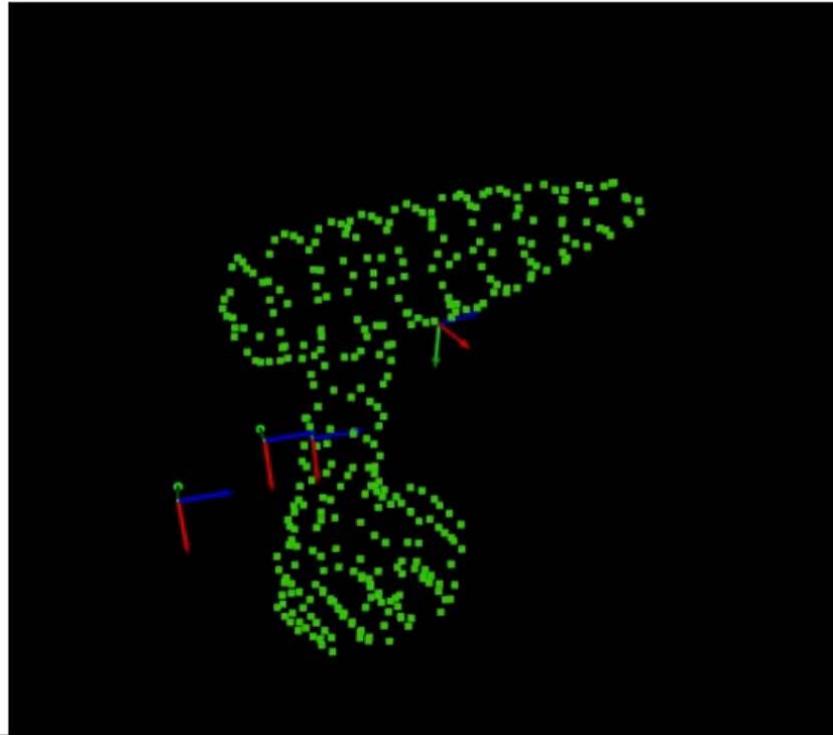
# Shape-aware Non-rigid Registration

- Partial view of novel instance
- Deformed canonical model



# Transference of Grasping Skills

Warp grasping information



# Grasping an Unknown Power Drill

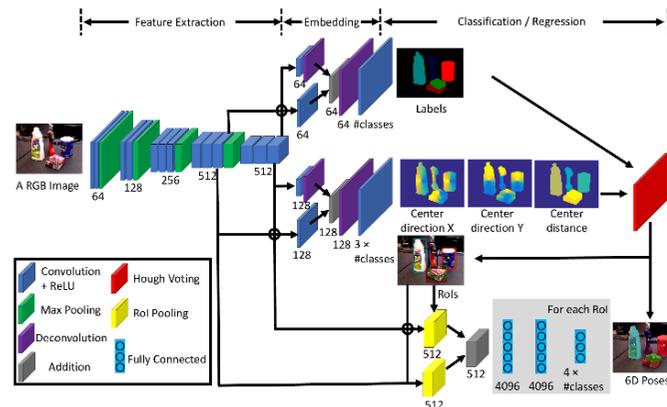


# Fully Convolutional 6D Pose Estimation

## Extending PoseCNN

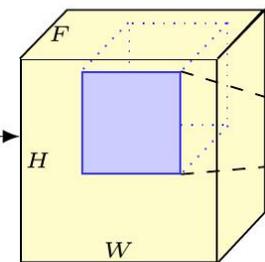
[Xiang, Schmidt, Narayanan, Fox: PoseCNN: A convolutional neural network for 6D object pose estimation in cluttered scenes. RSS 2018]

- Fully convolutional (per-pixel) prediction of pose parameters: 2D center offset, depth, orientation as quaternion

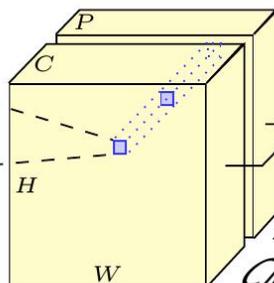


Input

Backbone network



Feature maps



Outputs

*Pose Class prob.*

Post-processing

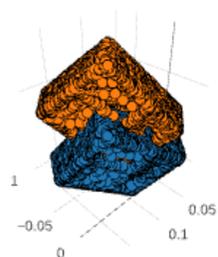
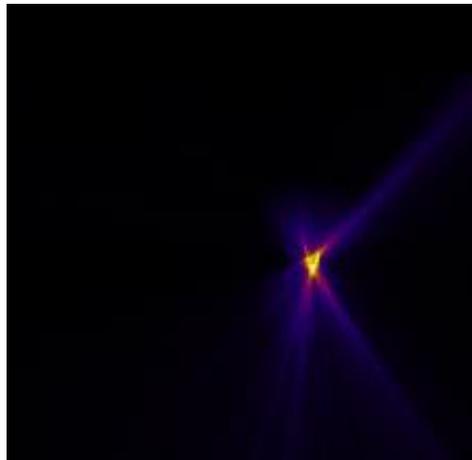
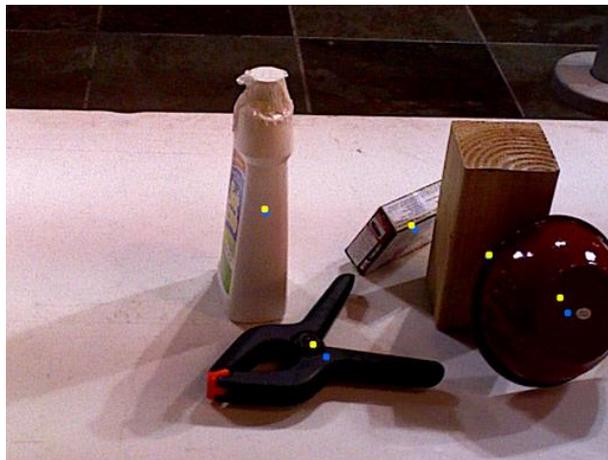


Object contours with poses

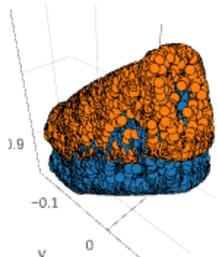
mouse\_traps  
conf: 0.916822  
windex  
conf: 0.902561  
q-tips\_500  
conf: 0.361778  
green\_binder  
conf: 0.310353  
green\_wine\_glass  
conf: 0.560962  
bubble\_microphone  
conf: 0.363494  
fiskars\_scissors  
conf: 0.838783  
band\_aid\_tape  
conf: 0.346222

# Fully Convolutional 6D Pose Estimation

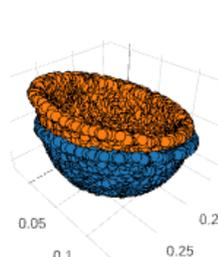
- Hough voting to find object centers in 2D
- Quaternions aggregated using Hough inliers and semantic segmentation



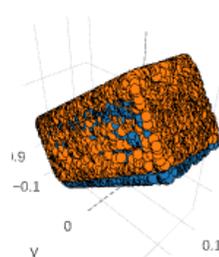
• ground truth  
• estimated



• ground truth  
• estimated



• ground truth  
• estimated



• ground truth  
• estimated



# From Turntable Captures to Meshes

## ■ Turntable setup:

- DSLR (Nikon D3400)
- Depth sensor (PrimeSense Carmine)
- Dynamixel actuator

## ■ Fast calibration:

- Automatic capture of checkerboard images
- Ceres-based optimization of camera extrinsics

## ■ Meshing:

- Masking via Background Subtraction
- Extract isosurface of visual hull + TSDF from depth sensor
- Texturing with Color Map Optimization



### Range Sensor and Silhouette Fusion for High-Quality 3D Scanning

Karthik S. Narayan, James Sha, Arjun Singh, and Pieter Abbeel  
ICRA 2015

### Color Map Optimization for 3D Reconstruction with Consumer Depth Cameras

Qian-Yi Zhou, Vladlen Koltun  
ACM TOG 2014

# From Turntable Captures to Meshes

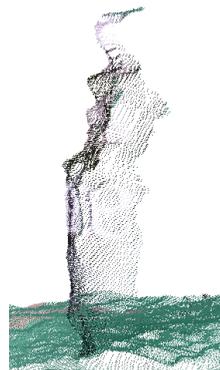


Fused & textured result

## Reflections & Unreliable Depth



DSLR image



Depth camera



Visual hull

## Concavities



Visual hull



Fused result

# An Alternative: CAD Meshes

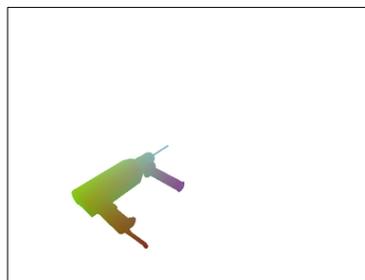
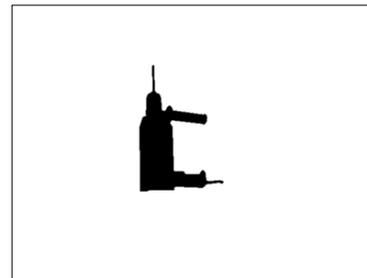
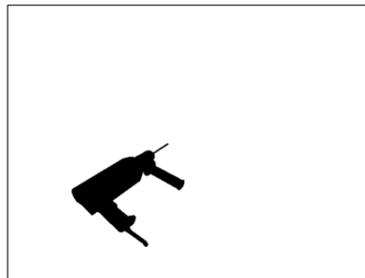
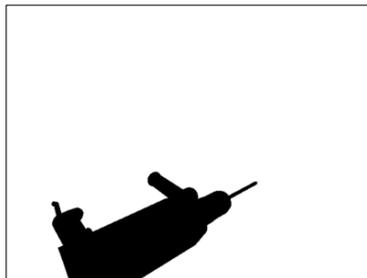
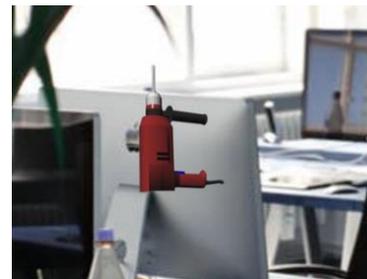
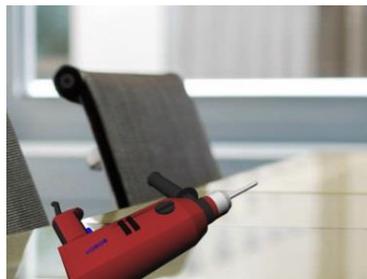
- Meshes found by search term "drill" on <https://sketchfab.com>



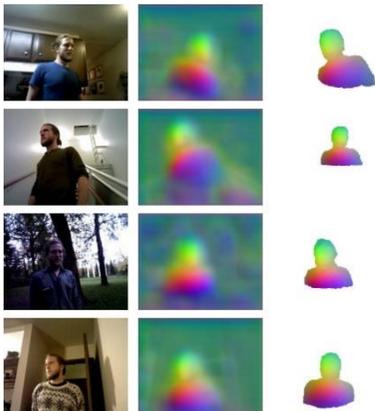
# Rendering 3D Scenes

## Advantages of mesh-based scene synthesis:

- Generate new scenes on-the-fly during training
- OpenGL/CUDA interop
- Simulate variations in hue, lighting, scale, rotation, camera intrinsics, ...
- Cheap ground truth:
  - segmentation labels
  - object-centric coordinates
  - occlusion information
  - ...



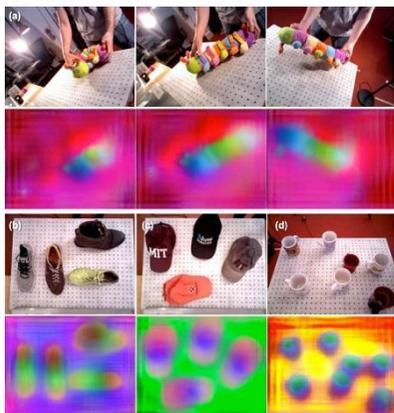
# Self-Supervised Feature Learning



Self-supervised Visual Descriptor Learning for Dense Correspondence

T. Schmidt, R. A. Newcombe, D. Fox

Robotics and Automation Letters, 2017



Dense Object Nets: Learning Dense Visual Object Descriptors

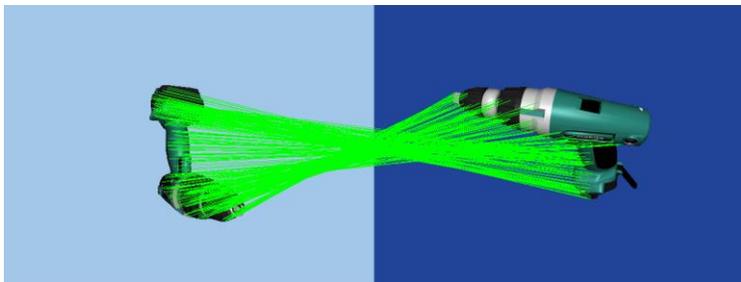
by and for Robotic Manipulation

P. Florence, L. Manuelli, R. Tedrake

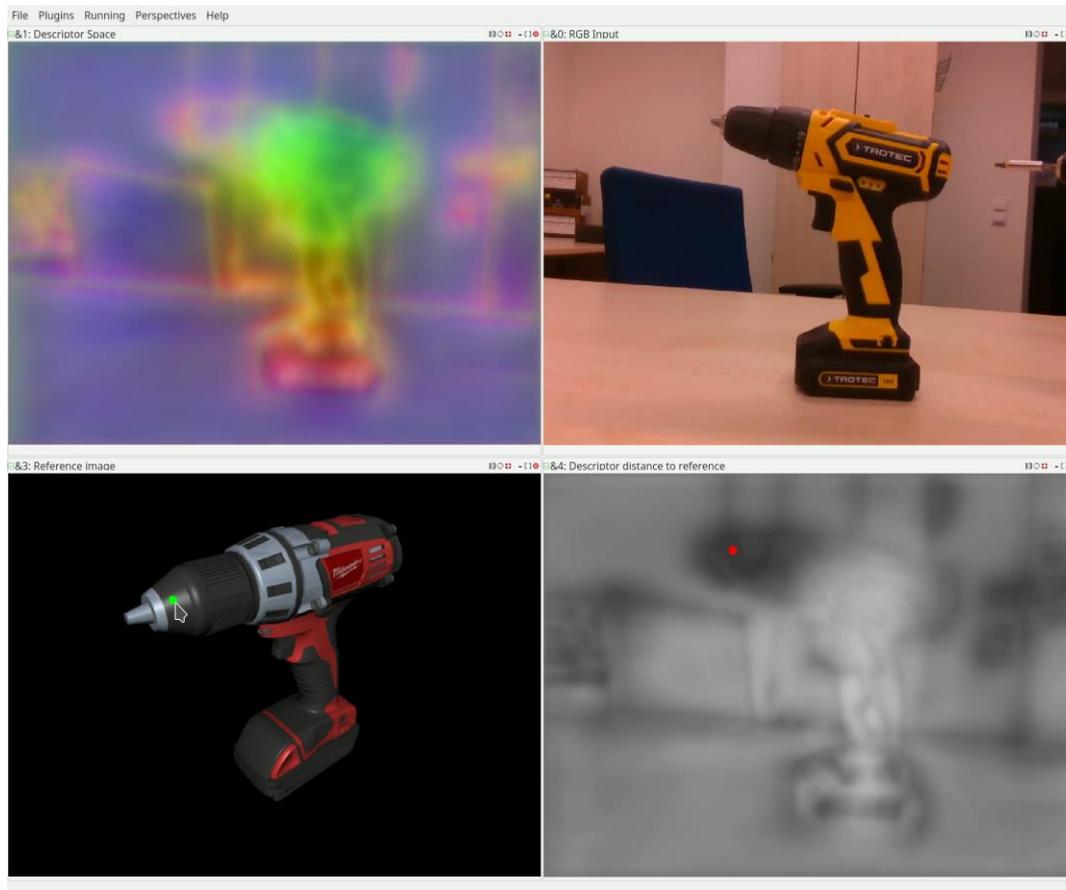
CoRL 2018

# Visual Descriptor Learning

- Trained on 1100 frames rendered from 22 CAD meshes

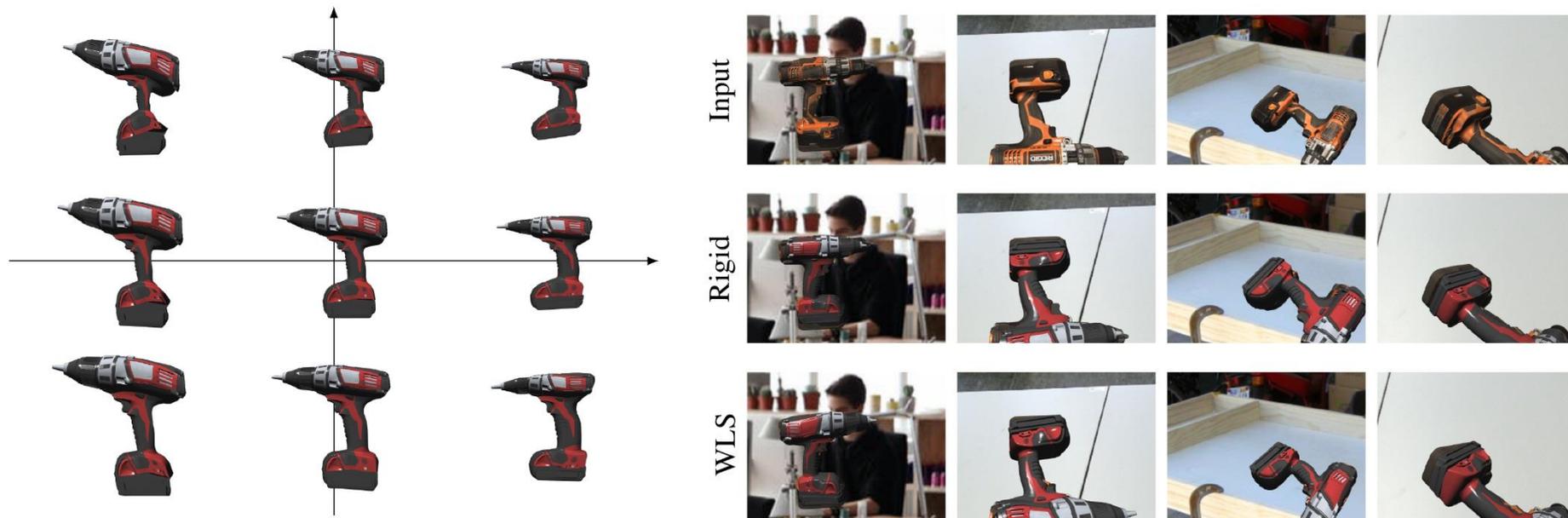


- Pixel-wise contrastive loss
- No training signal between different instances!



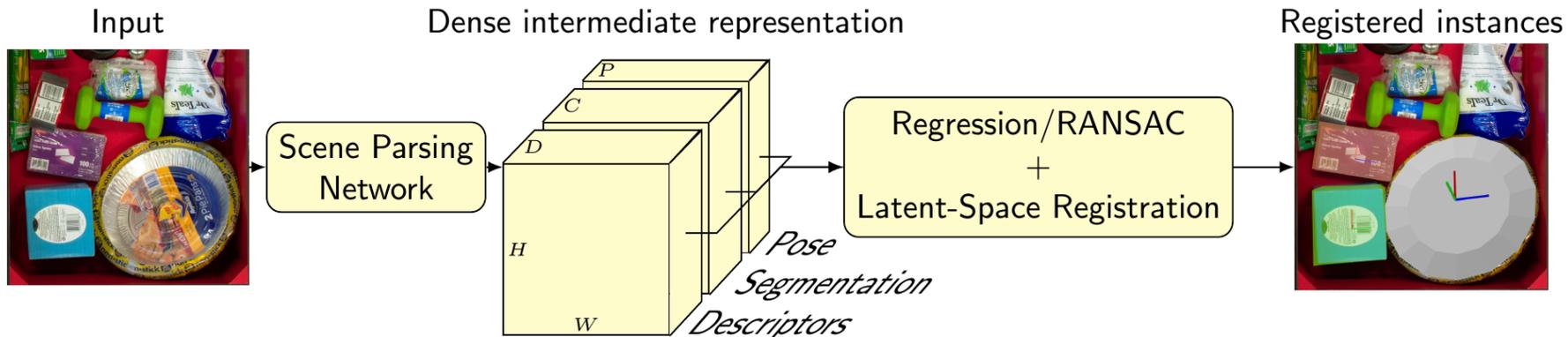
# Combination with Non-Rigid Registration

- Shape space creation using correspondence from visual descriptors
- Inference: Semantic segmentation, RANSAC, shape-aware fine registration

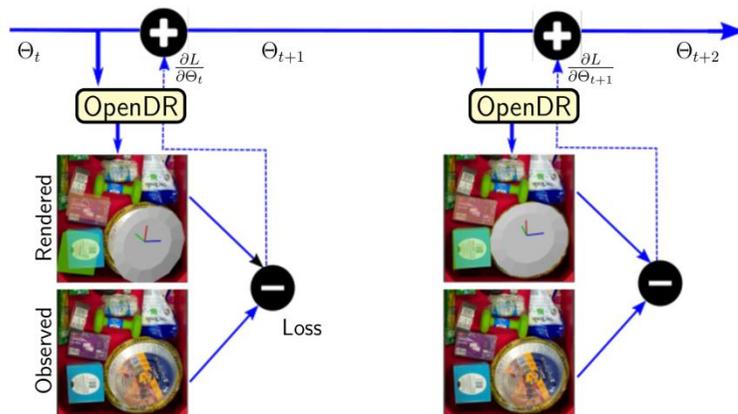


# Outlook

- Complete scene parsing pipeline utilizing learned descriptors and shape models

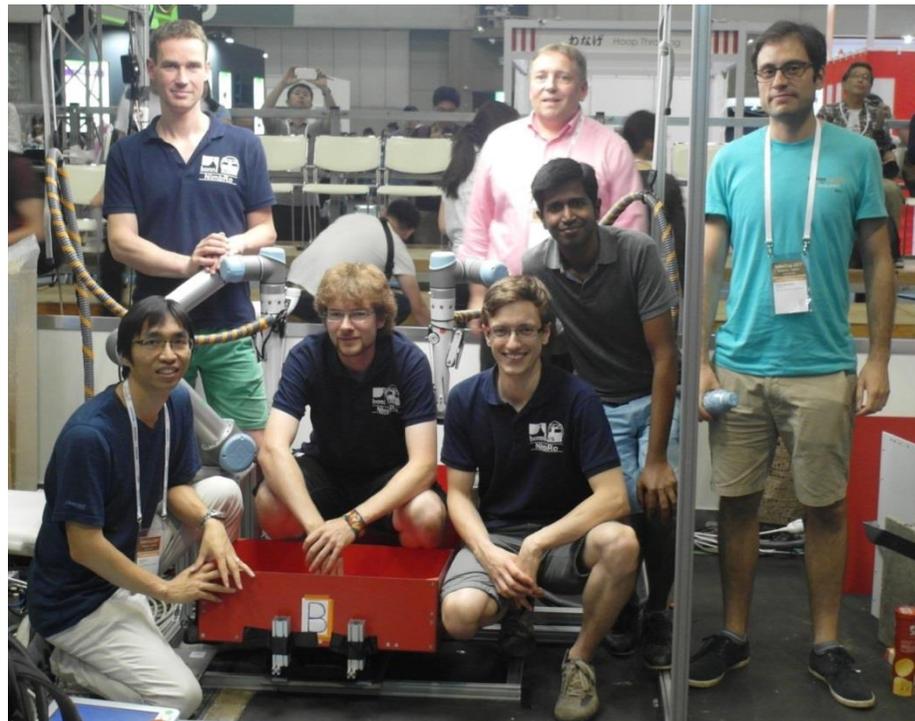


- Iterative refinement by render and compare



# Conclusions

- Developed methods for learning semantic perception of cluttered bin scenes
  - Data capture and annotation
  - Synthesizing scenes
  - Deformable models
- Integration to APC/ARC systems
- Addressed data problem by
  - Data capture and annotation
  - Synthesizing scenes
  - Deformable models
- Much further research needed for complete scene understanding



ARC 2017 team Nimbro Picking