# Learning Semantic Environment Perception for Cognitive Robots

**Sven Behnke**

University of Bonn, Germany

Computer Science Institute VI

Autonomous Intelligent Systems

# Some of Our Cognitive Robots

- Equipped with many sensors and DoFs
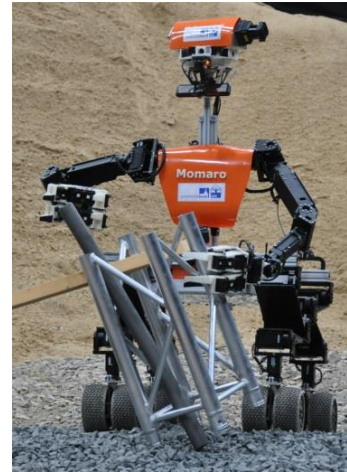- Demonstration in complex scenarios
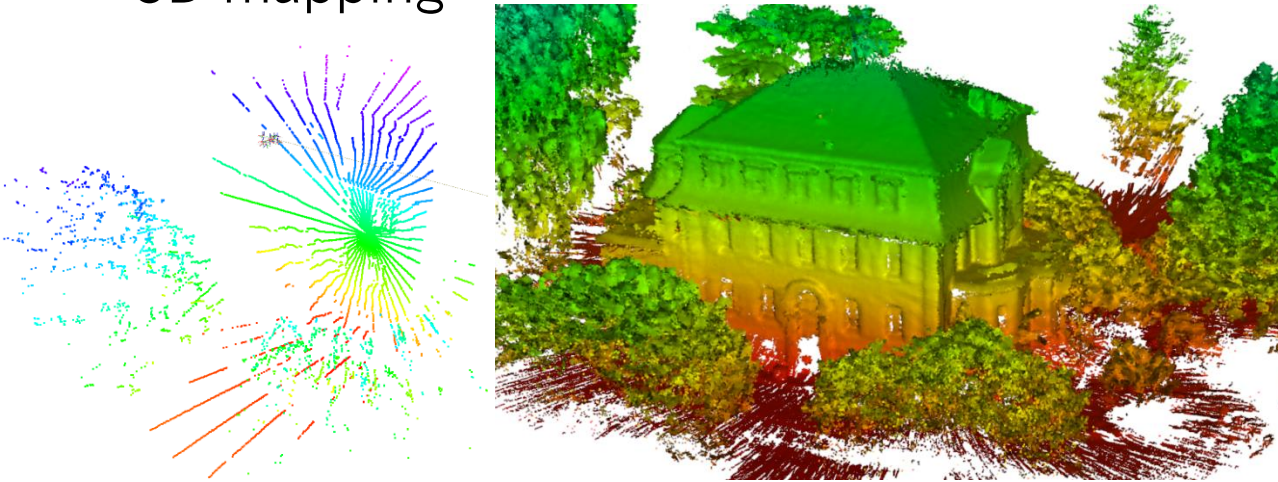


MAV



Soccer robot



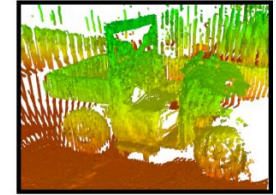Service robot



Exploration robot



Picking robot

Sven Behnke: Semantic Environment Perception

universität**bonn** **ais**

# 3D Environment Perception

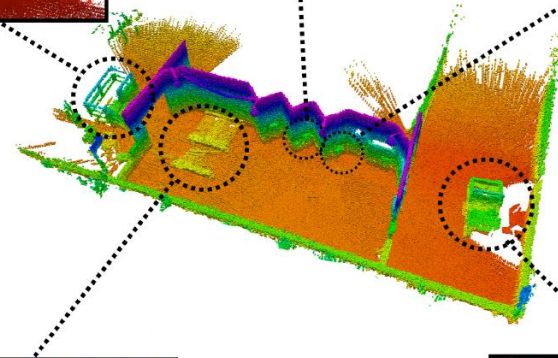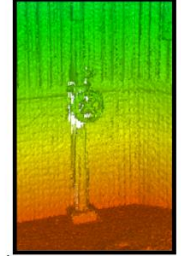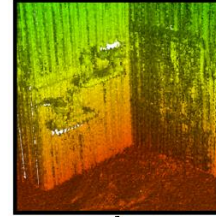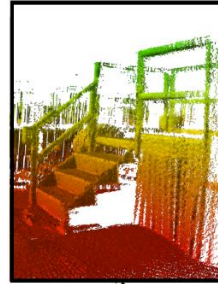- 3D laser scanner, dual wide-angle stereo cameras, ultrasound, Quad Core i7
- Autonomous navigation close to structures
- 3D mapping



[Droeschel et al. JFR 2016]

universität**bonn** **ais**

# 3D Mapping

- Registering 3D laser scans



[Droeschel et al. 2016]

Sven Behnke: Semantic Environment Perception

universität**bonn** **ais**

# Mobile Manipulation in Mars-like Environment



Sven Behnke: Semantic Environment Perception

# Autonomous Mission Execution

- 3D mapping, localization, mission and navigation planning



- 3D object perception and grasping



Sven Behnke: Semantic Environment Perception

[Schwarz et al. Frontiers 2016]

universität**bonn** ais

# Cognitive Service Robot Cosero



Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Table-top Analysis and Grasp Planning
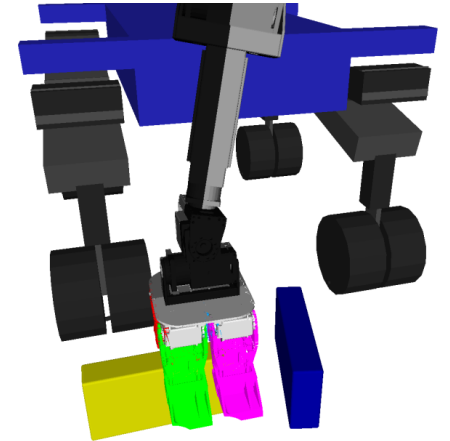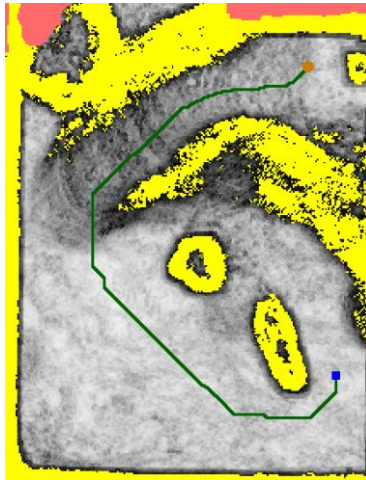
- Detection of clusters above horizontal plane

- Two grasps (top, side)



- Flexible grasping of many unknown objects



[Stückler et al, Robotics and Autonomous Systems, 2013]

Sven Behnke: Semantic Environment Perception

universität**bonn** **ais**

# 3D Mapping by RGB-D SLAM

- Modelling of shape and color distributions in voxels
- Local multiresolution
- Efficient registration of views on CPU
- Global optimization

- Multi-camera SLAM

5cm

2,5cm

[Stoucken]

universität**bonn** ais

# Learning and Tracking Object Models

- Modeling of objects by RGB-D-SLAM



- Real-time registration with current RGB-D frame



Sven Behnke: Semantic Environment Perception

# Deformable RGB-D-Registration

- Based on Coherent Point Drift method [Myronenko & Song, PAMI 2010]
- Multiresolution Surfel Map allows real-time registration



Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Transformation of Poses on Object

- Derived from the deformation field



[Stückler, Behnke, ICRA2014]

universität**bonn** ais

# Grasp & Motion Skill Transfer



[Stückler, Behnke, ICRA2014]

Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Tool use: Bottle Opener

- Tool tip perception



- Extension of arm kinematics
- Perception of crown cap
- Motion adaptation



[Stückler, Behnke, Humanoids 2014]

Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Hierarchical Object Discovery trough Motion Segmentation

- Simultaneous object modeling and motion segmentation



- Inference of a segment hierarchy



[Stückler, Behnke: IJCAI 2013]

Sven Behnke: Semantic Environment Perception

# Semantic Mapping

- Pixel-wise classification of RGB-D images by random forests
- Compare color / depth of regions
- Size normalization
- 3D fusion through RGB-D SLAM
- Evaluation on NYU depth v2



$$p(c|\mathcal{F}, q) = \frac{1}{K} \sum_{k=1}^{K} p(c|l_k(q))$$

[Stückler, Biresev, Behnke: IROS 2012]

Ground truth

Segmentation

| Accuracy in % | Ø Classes | Ø Pixels |
|---|---|---|
| Silberman et al. 2012 | 59,6 | 58,6 |
| Couprie et al. 2013 | 63,5 | 64,5 |
| Random forest | 65,0 | 68,1 |
| 3D-Fusion | **66,8** | |

Sven Behnke: Semantic Environment Perception

universität**bonn**   ais

# Learning Depth-sensitive CRFs

- SLIC+depth super pixels
- Unary features: random forest
- Height feature



- Pairwise features
  - Color contrast
  - Vertical alignment
  - Depth difference
  - Normal differences



- Results:

|  | class average | pixel average |
|---|---|---|
| RF | 65.0 | 68.3 |
| RF + SP | 65.7 | 70.1 |
| RF + SP + SVM | 70.4 | 70.3 |
| RF + SP + CRF | **71.9** | **72.3** |
| Silberman *et al.* | 59.6 | 58.6 |
| Couprie *et al.* | 63.5 | 64.5 |



RGB Image

Depth Image

Random Forest

Superpixel

3D Point Cloud

Unary Features

Pairwise Features

CRF Prediction



Random forest

CRF prediction

Ground truth

Sven Behnke: Semantic Environment Perception

universität**bonn**  ais

# Deep Learning

- Learning layered represen- tations



[Schulz; Behnke, KI 2012]

Sven Behnke: Semantic Environment Perception

# Object-class Segmentation

- Class annotation per pixel

- Multi-scale input channels

- Evaluated on MSRC-9/21 and INRIA Graz-02 data sets



Input    Output    Truth

Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Object Detection in Natural Images

- Bounding box annotation
- Structured loss that directly maximizes overlap of the prediction with ground truth bounding boxes
- Evaluated on two of the Pascal VOC 2007 classes



[Schulz, Behnke, ICANN 2014]

Sven Behnke: Semantic Environment Perception

# RGB-D Object-Class Segmentation

- Covering windows segmented with CNN
- Scale input according to depth, compute pixel hight



RGB   Depth   Height   Truth   Output

| Method | floor | struct | furnit | prop | Class Avg. | Pixel Acc. |
|---|---|---|---|---|---|---|
| CW | 84.6 | 70.3 | 58.7 | 52.9 | 66.6 | 65.4 |
| CW+DN | 87.7 | 70.8 | 57.0 | 53.6 | 67.3 | 65.5 |
| CW+H | 78.4 | 74.5 | 55.6 | 62.7 | 67.8 | 66.5 |
| CW+DN+H | 93.7 | 72.5 | 61.7 | 55.5 | 70.9 | 70.5 |
| CW+DN+H+SP | 91.8 | 74.1 | 59.4 | 63.4 | 72.2 | 71.9 |
| CW+DN+H+CRF | 93.5 | 80.2 | 66.4 | 54.9 | **73.7** | **73.4** |
| Müller et al.[8] | 94.9 | 78.9 | 71.1 | 42.7 | 71.9 | 72.3 |
| Random Forest [8] | 90.8 | 81.6 | 67.9 | 19.9 | 65.1 | 68.3 |
| Couprie et al.[9] | 87.3 | 86.1 | 45.5 | 35.5 | 63.6 | 64.5 |
| Höft et al.[10] | 77.9 | 65.4 | 55.9 | 49.9 | 62.3 | 62.0 |
| Silberman [12] | 68 | 59 | 70 | 42 | 59.7 | 58.6 |

CW is covering windows, H is height above ground, DN is depth normalized patch sizes. SP is averaged within superpixels and SVM-reweighted. CRF is a conditional random field over superpixels [8]. Structure class numbers are optimized for class accuracy.

[Schulz, Höft, Behnke, ESANN 2015]

universität**bonn** **ais**

# Neural Abstraction Pyramid



[Behnke, Rojas, IJCNN 1998]
[Behnke, LNCS 2766, 2003]

Abstract features

- Data-driven
- Analysis
- Feature extraction

- Model-driven
- Synthesis
- Feature expansion

Signals

- Grouping  - Competition  - Completion

Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Iterative Image Interpretation

- Interpret most obvious parts first
- Use partial interpretation as context to resolve local ambiguities

Sven Behnke: Semantic Environment Perception

# Neural Abstraction Pyramid for RGB-D Video Object-class Segmentation

- Recursive computation is efficient for temporal integration



Neural Abstraction Pyramid

[Pavel, Schulz, Behnke, Neural Networks 2017]

universität**bonn** ais

# Geometric and Semantic Features for RGB-D Object-class Segmentation

- New **geometric** feature: distance from wall
- **Semantic** features pretrained from ImageNet
- Both help significantly

[Husain et al. RA-L 2016]



RGB    Truth    DistWall    OutWO    OutWithDistWall

Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Semantic Segmentation Priors for Object Discovery

- Combine bottom-up object discovery and semantic priors
- Semantic segmentation used to classify color and depth superpixels
- Higher recall, more precise object borders

[Garcia et al. ICPR 2016]

universität**bonn** ais

# RGB-D Object Recognition and Pose Estimation



[Schwarz, Schulz, Behnke, ICRA2015]

Sven Behnke: Semantic Environment Perception

# Canonical View, Colorization

- Objects viewed from different elevation
- Render canonical view

- Colorization based on distance from center vertical



[Schwarz, Schulz, Behnke, ICRA2015]

Sven Behnke: Semantic Environment Perception

universität**bonn** **ais**

# Pretrained Features Disentangle Data

- t-SNE embedding



[Schwarz, Schulz, Behnke ICRA2015]

Sven Behnke: Semantic Environment Perception

# Recognition Accuracy

- Improved both category and instance recognition

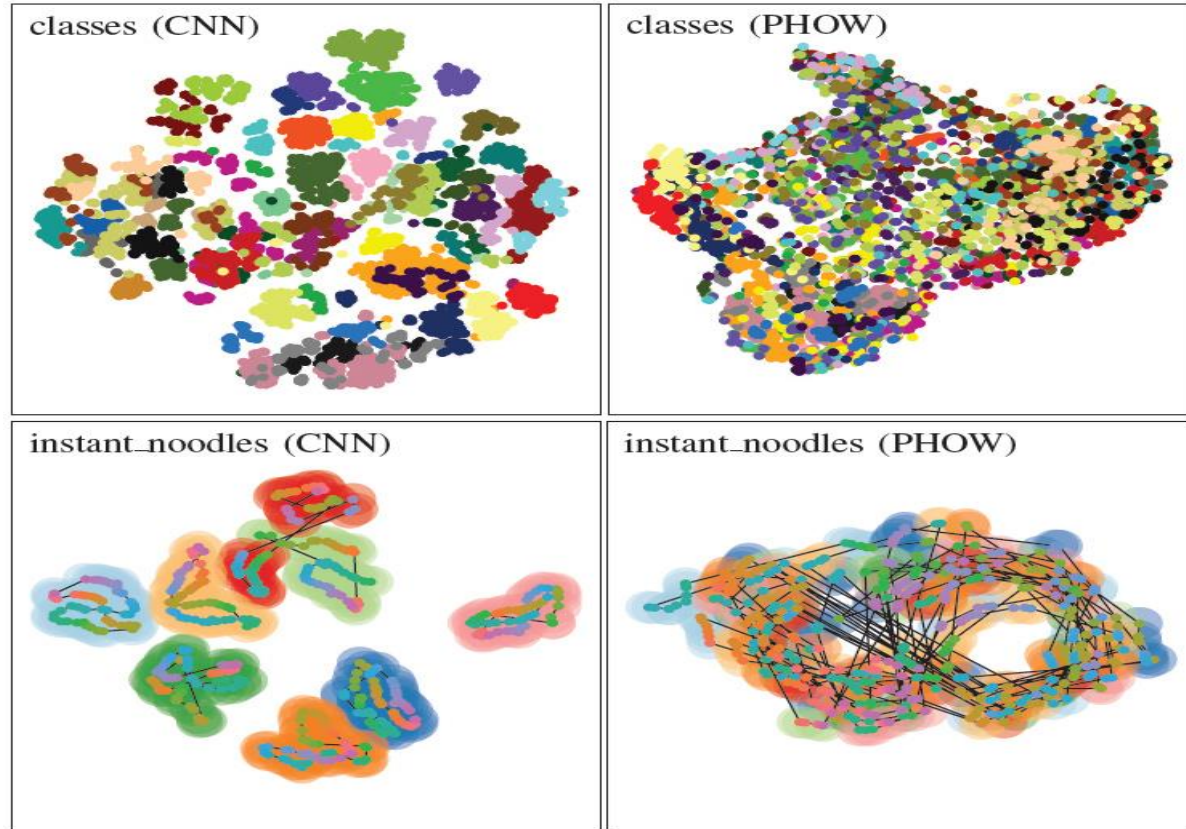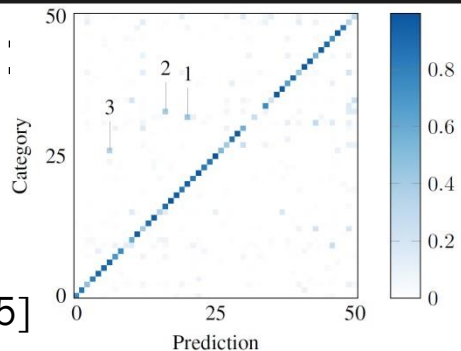| Method | Category Accuracy (%) | | Instance Accuracy (%) | |
|---|---|---|---|---|
| | RGB | RGB-D | RGB | RGB-D |
| Lai *et al.* [1] | $74.3 \pm 3.3$ | $81.9 \pm 2.8$ | 59.3 | 73.9 |
| Bo *et al.* [2] | $82.4 \pm 3.1$ | $87.5 \pm 2.9$ | **92.1** | 92.8 |
| PHOW[3] | $80.2 \pm 1.8$ | — | 62.8 | — |
| **Ours** | $\mathbf{83.1 \pm 2.0}$ | $88.3 \pm 1.5$ | 92.0 | **94.1** |
| **Ours** | $\mathbf{83.1 \pm 2.0}$ | $\mathbf{89.4 \pm 1.3}$ | 92.0 | **94.1** |

- Confusion:



1: pitcher / coffe mug



2: peach / sponge



[Schwarz, Schulz, Behnke, ICRA2015]

Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Amazon Picking Challenge

- Large variety of objects
- Unordered in shelf or tote
- Picking and stowing tasks



[Schwarz et al. ICRA 2017]

Sven Behnke: Semantic Environment Perception

# Deep Learning Semantic Segmentation

- Adapted from our segmentation of indoor scenes [Husain et al. RA-L 2016]



[Schwarz et al. ICRA 2017]

Sven Behnke: Semantic Environment Perception

# DenseCap Object Detection



[Schwarz et al. ICRA 2017]　　　　　　　　　　　　　[Johnson et a CVPR 2016]

Sven Behnke: Semantic Environment Perception

# Combined Detection and Segmentation

Detection



x →

[Schwarz et al. IJRR 2017]

Segmentation

Sven Behnke: Semantic Environment Perception

universität**bonn** | ais

# Stowing



Sven Behnke: Semantic Environment Perception

# Picking



Sven Behnke: Semantic Environment Perception

# NimbRo Picking APC 2016 Results

- 2$^{nd}$ Place Stowing (186 points)
- 3$^{rd}$ Place Picking (97 points)

[Schwarz et al. IJRR 2017]



Sven Behnke: Semantic Environment Perception

# Detection of Tools



Sven Behnke: Semantic Environment Perception [Schwarz et al. IJRR 2017]

universität**bonn** ais

# MBZIRC Challenge 2



Sven Behnke: Semantic Environment Perception

# Wrench Selection: Detection of Tool Ends



Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Amazon Robotics Challenge 2017

- Training with rendered scenes



Sven Behnke: Semantic Environment Perception

# **Conclusions**

- Semantic perception is challenging
- Simple methods rely on strong assumptions
- Depth helps with segmentation, allows for size normalization, geometric features, shape descriptors
- Deep learning methods work well
- Transfer of features from large data sets
- Synthetic training
- Many open problems, e.g. total scene understanding, incorporating physics, …

Sven Behnke: Semantic Environment Perception

universität**bonn** ais

# Questions?

Sven Behnke: Semantic Environment Perception