

# Perception, Planning, and Learning for Cognitive Robots

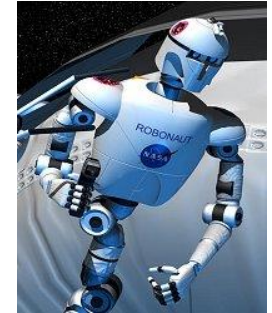
**Sven Behnke**

University of Bonn  
Computer Science Institute VI  
Autonomous Intelligent Systems



# Many New Application Areas for Robots

- Self-driving cars
- Logistics
- Agriculture, mining
- Collaborative production
- Personal assistance
- Space, search & rescue
- Healthcare
- Toys



**Need more cognitive abilities!**

# Some of our Cognitive Robots

- Equipped with numerous sensors and actuators
- Complex demonstration scenarios



Soccer



Domestic service



Mobile manipulation



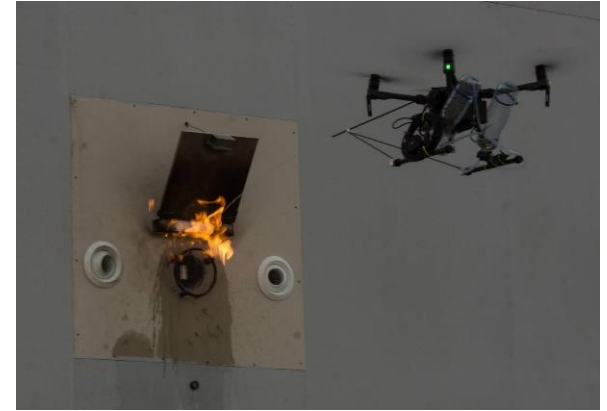
Bin picking



Aerial inspection

# Some more of our Cognitive Robots

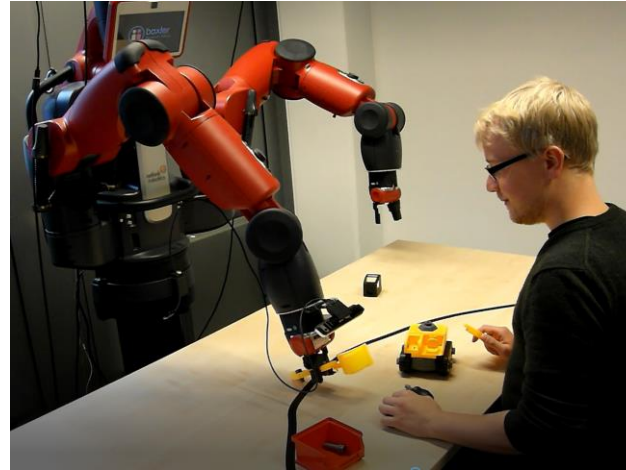
- Equipped with numerous sensors and actuators
- Complex demonstration scenarios



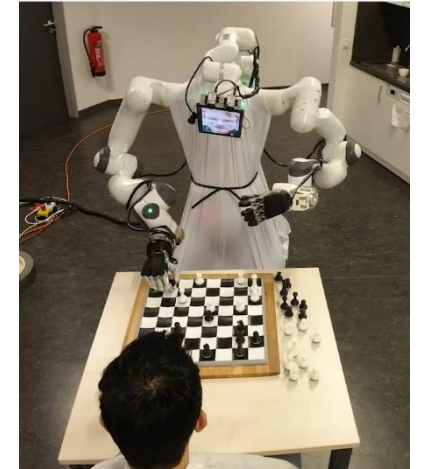
Rescue



Phenotyping



Human-robot collaboration

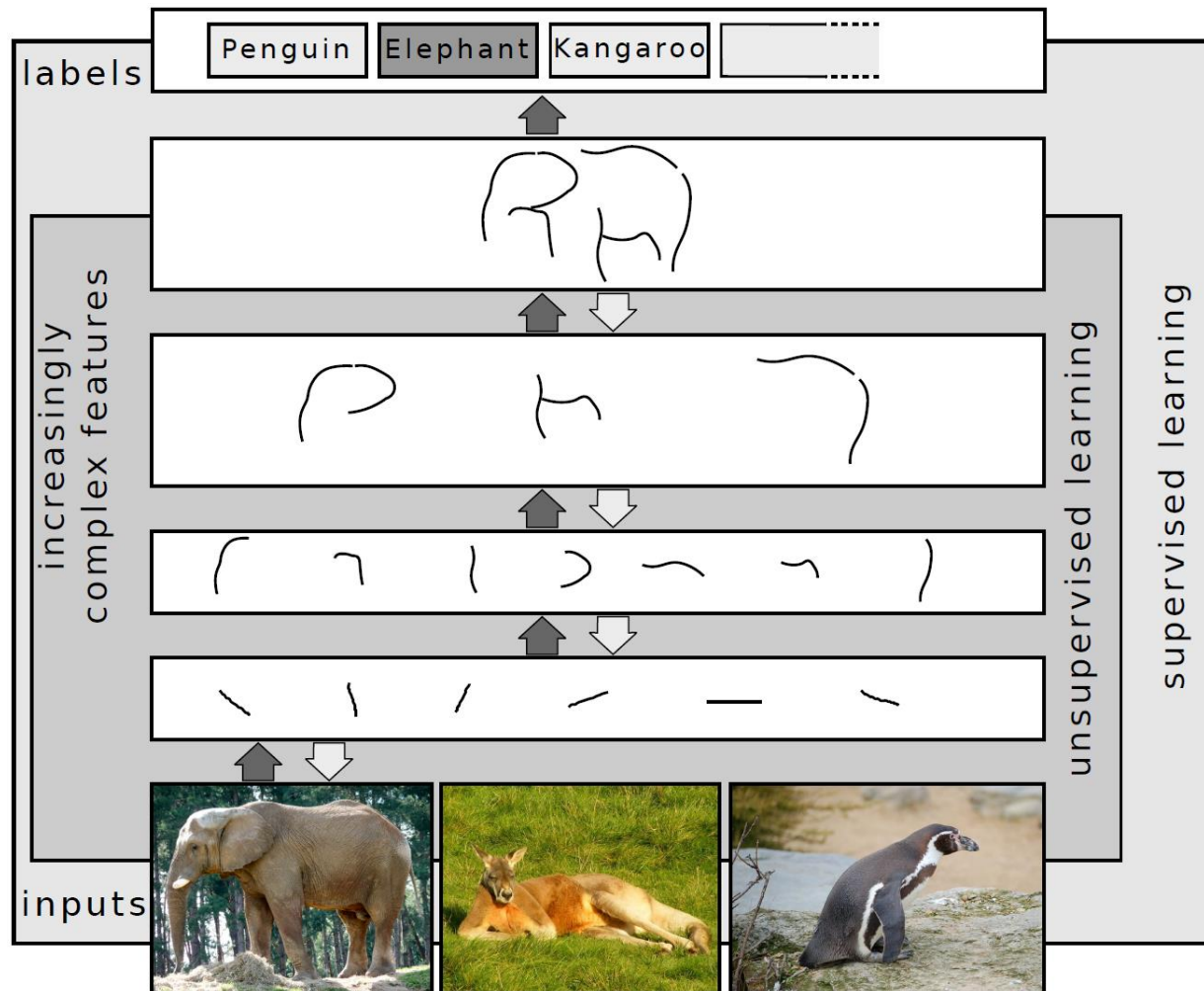


Telepresence

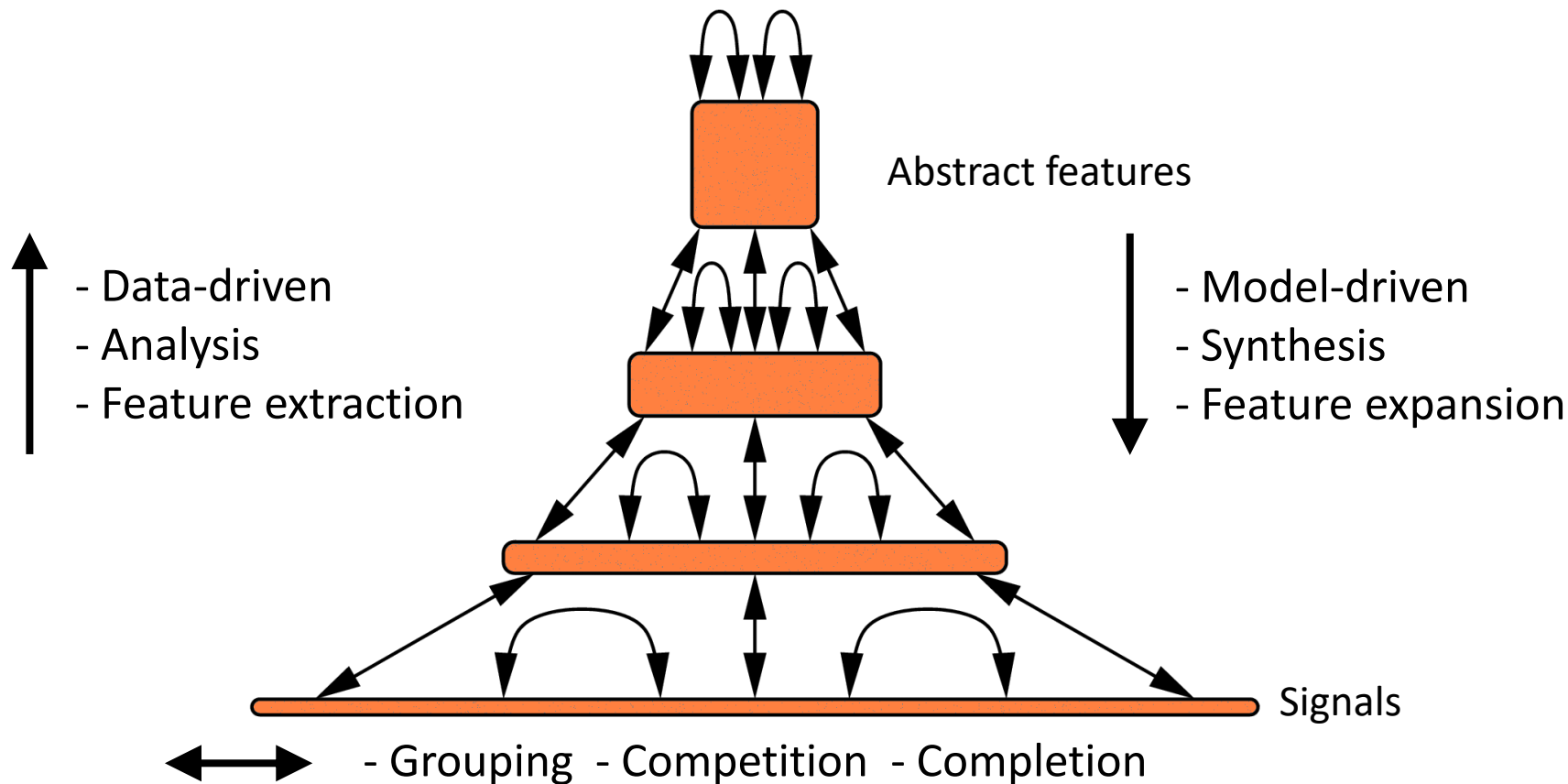
# Deep Learning

- Learning layered representations
- Compositionality

[Schulz;  
Behnke,  
KI 2012]

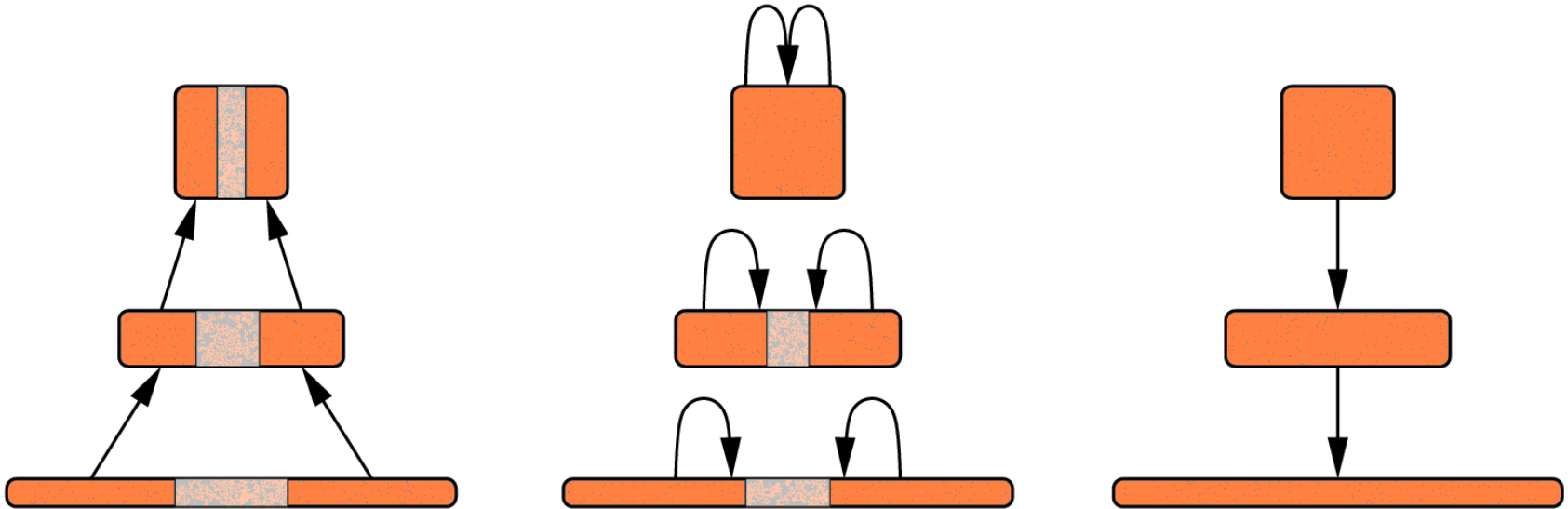


# Neural Abstraction Pyramid



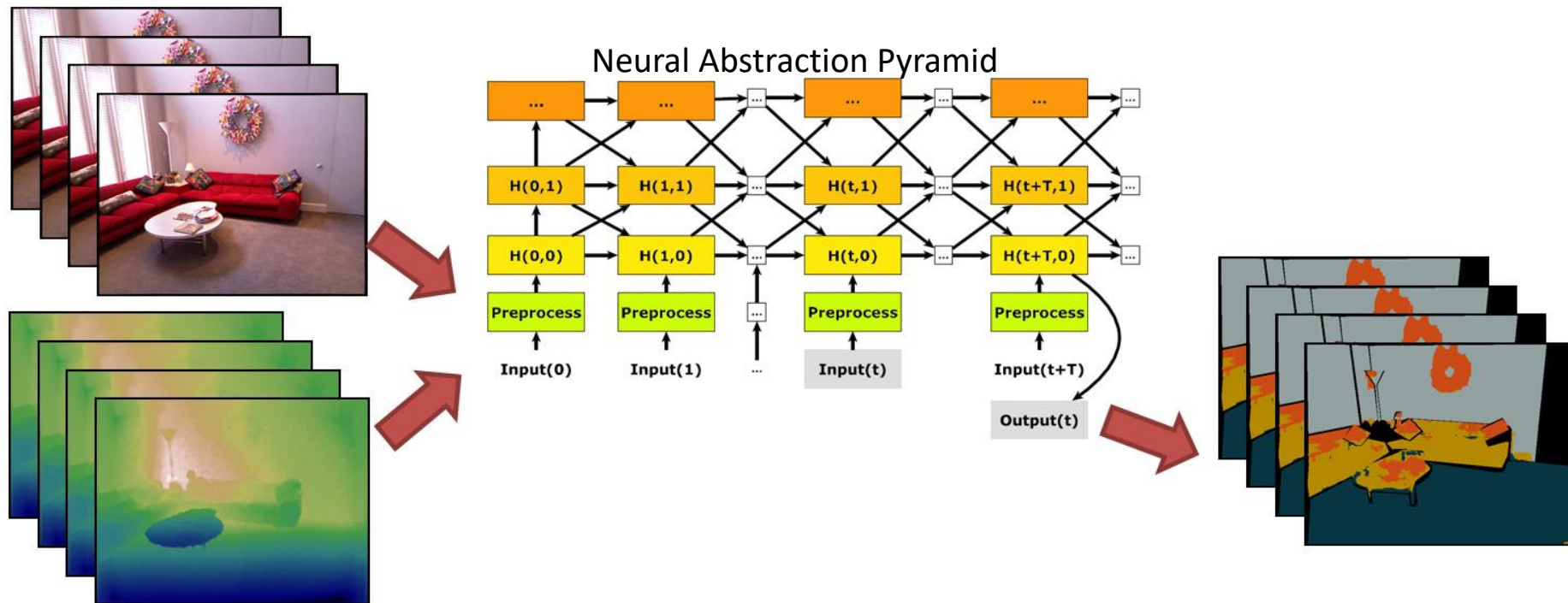
# Iterative Image Interpretation

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



# Neural Abstraction Pyramid for Semantic Segmentation of RGB-D Video

- Recursive computation is efficient for temporal integration





# The Data Problem

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

## 1. Transfer learning:

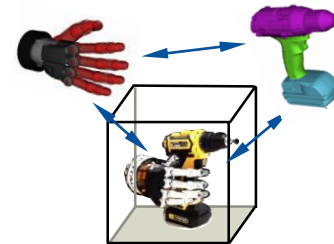
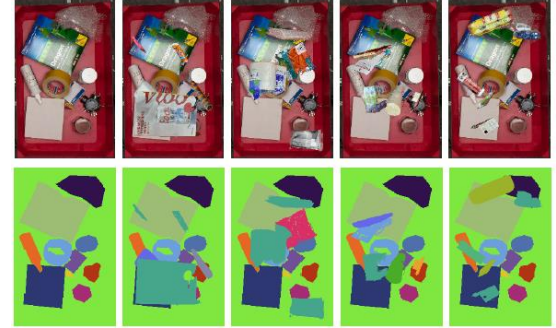
Pre-training on large related data,  
self-supervised learning

## 2. Generating data:

Online mesh databases,  
scene synthesis

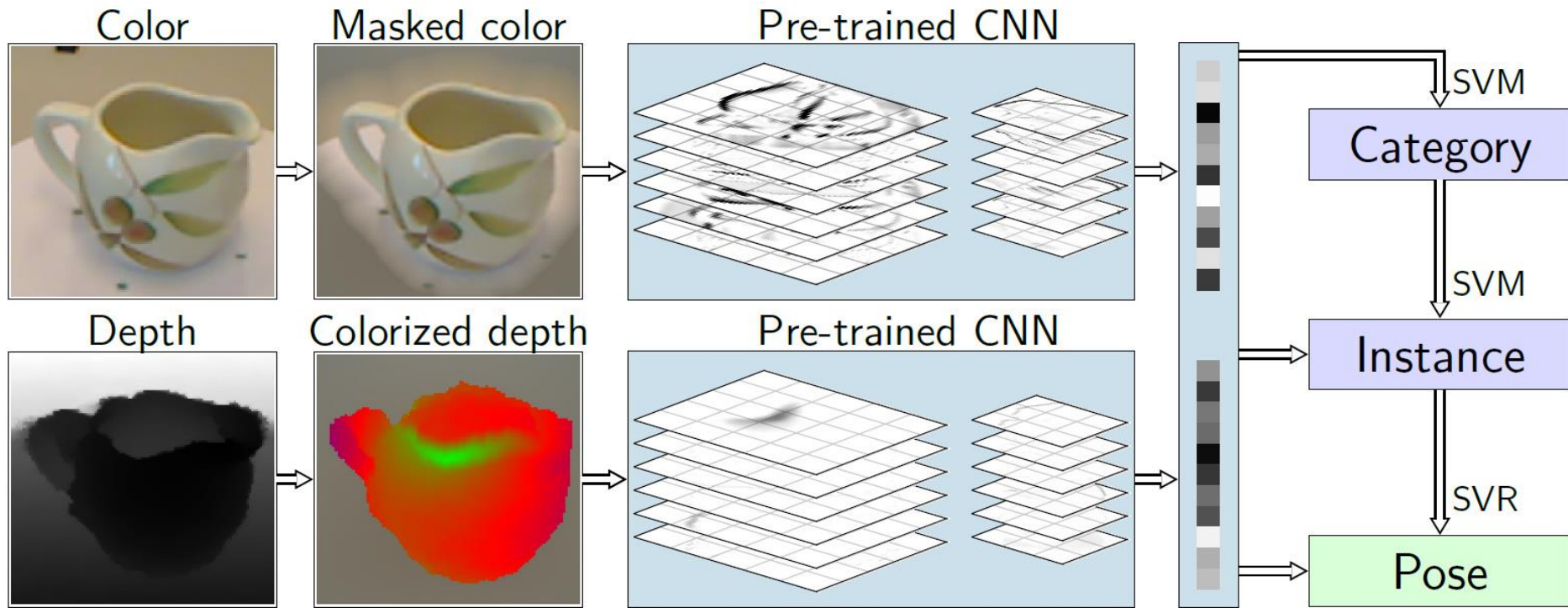
## 3. Inductive biases:

3D projective geometry,  
camera motion, canonical frames,  
object relations, compositionality, ...



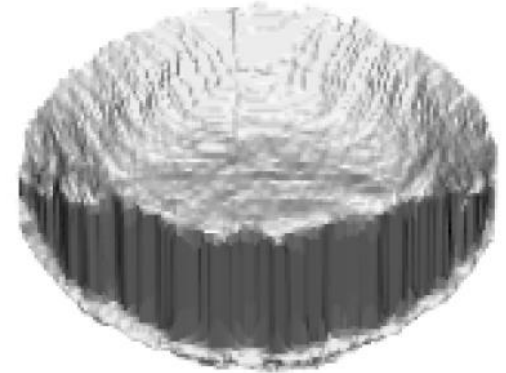
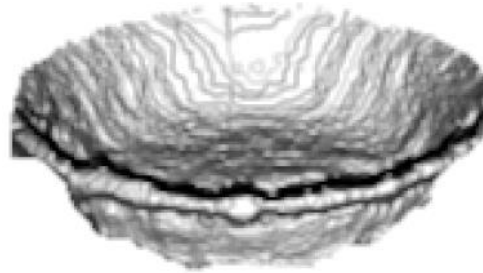
# RGB-D Object Recognition and Pose Estimation

- Transfer learning from large-scale data sets

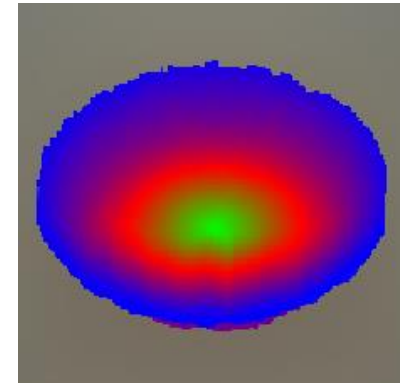
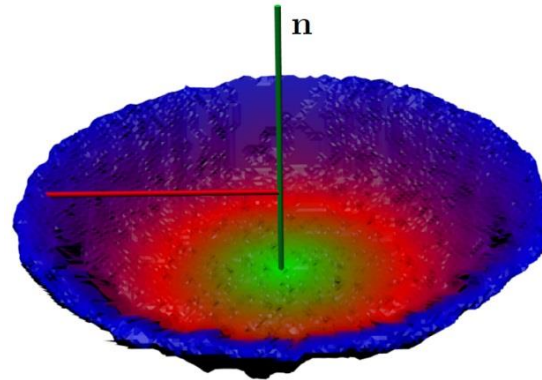


# Canonical View, Colorization

- Objects viewed from different elevation
- Render canonical view

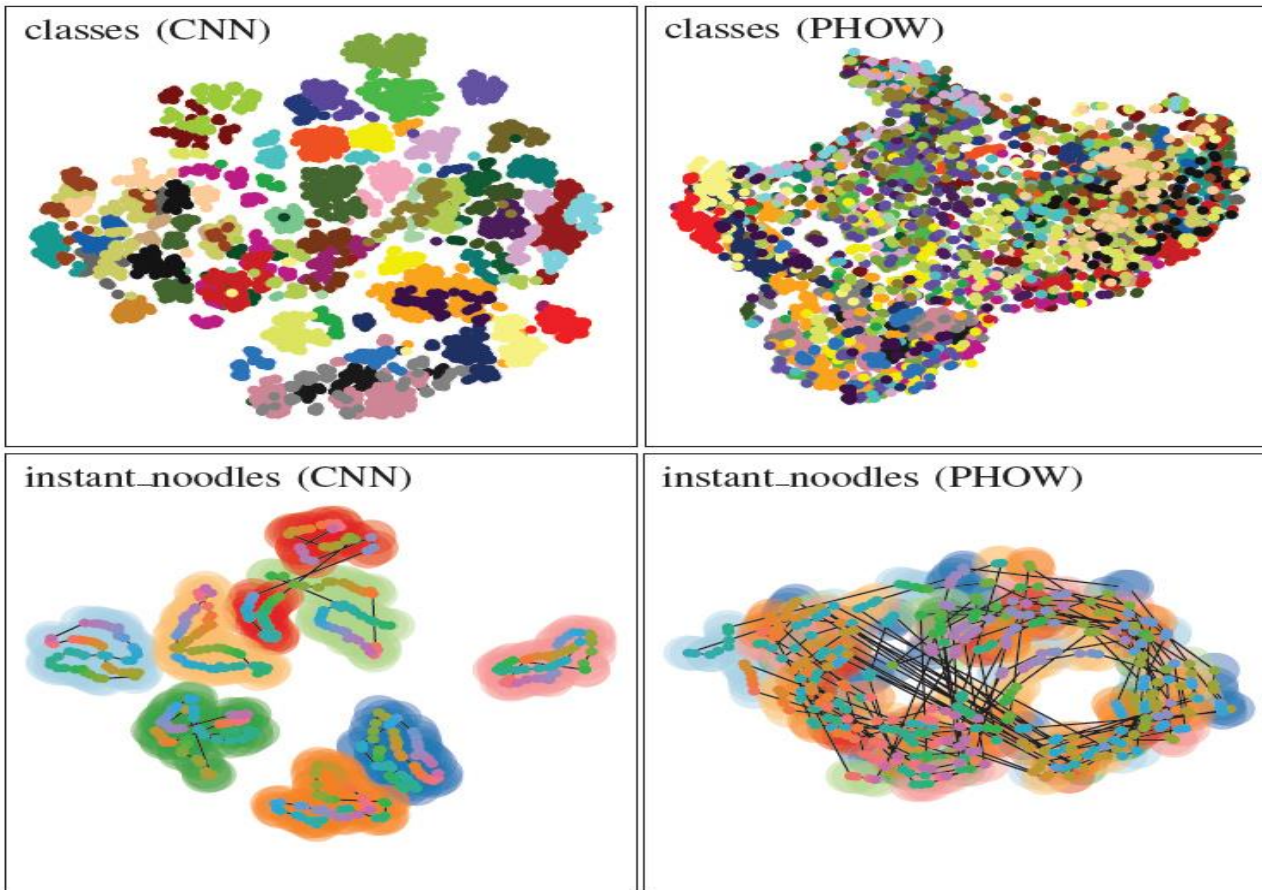


- Colorization based on distance from center vertical



# Pretrained Features Disentangle Data

- t-SNE embedding



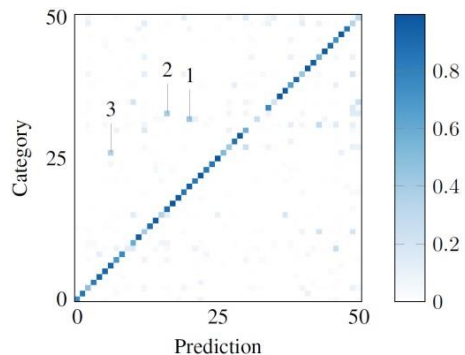
[Schwarz, Schulz,  
Behnke ICRA2015]

# Recognition Accuracy

## ■ Improved both category and instance recognition

Method	Category Accuracy (%)		Instance Accuracy (%)	
	RGB	RGB-D	RGB	RGB-D
Lai <i>et al.</i> [1]	74.3 ± 3.3	81.9 ± 2.8	59.3	73.9
Bo <i>et al.</i> [2]	82.4 ± 3.1	87.5 ± 2.9	<b>92.1</b>	92.8
PHOW[3]	80.2 ± 1.8	—	62.8	—
<b>Ours</b>	<b>83.1 ± 2.0</b>	88.3 ± 1.5	92.0	<b>94.1</b>
<b>Ours</b>	<b>83.1 ± 2.0</b>	<b>89.4 ± 1.3</b>	92.0	<b>94.1</b>

## ■ Confusion:



1: pitcher / coffee mug

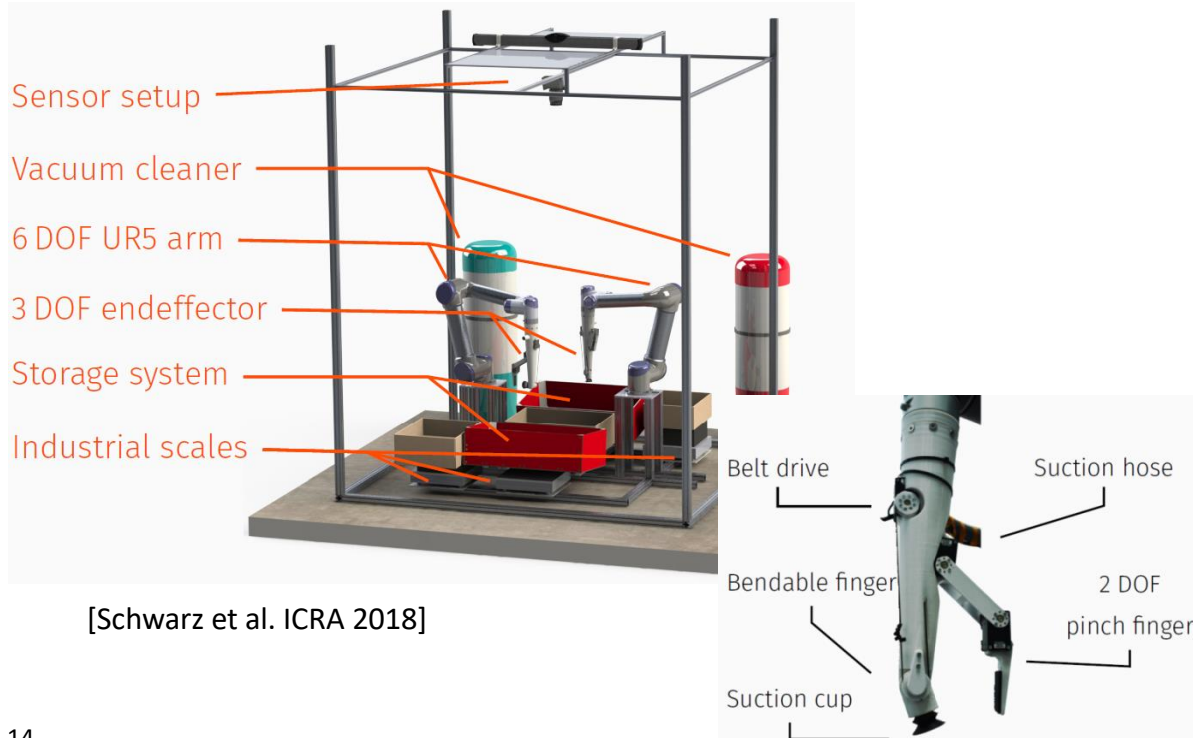


2: peach / sponge



# Amazon Robotics Challenge

- Storing and picking of items
- Dual-arm robotic system



[Amazon]

# Object Capture and Scene Rendering

## ■ Turntable + DSLR camera

## ■ Insertion in complex annotated scenes



# Semantic Segmentation and Grasp Pose Estimation

- Semantic segmentation using RefineNet [Lin et al. CVPR 2017]
- Grasp positions in segment centers



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

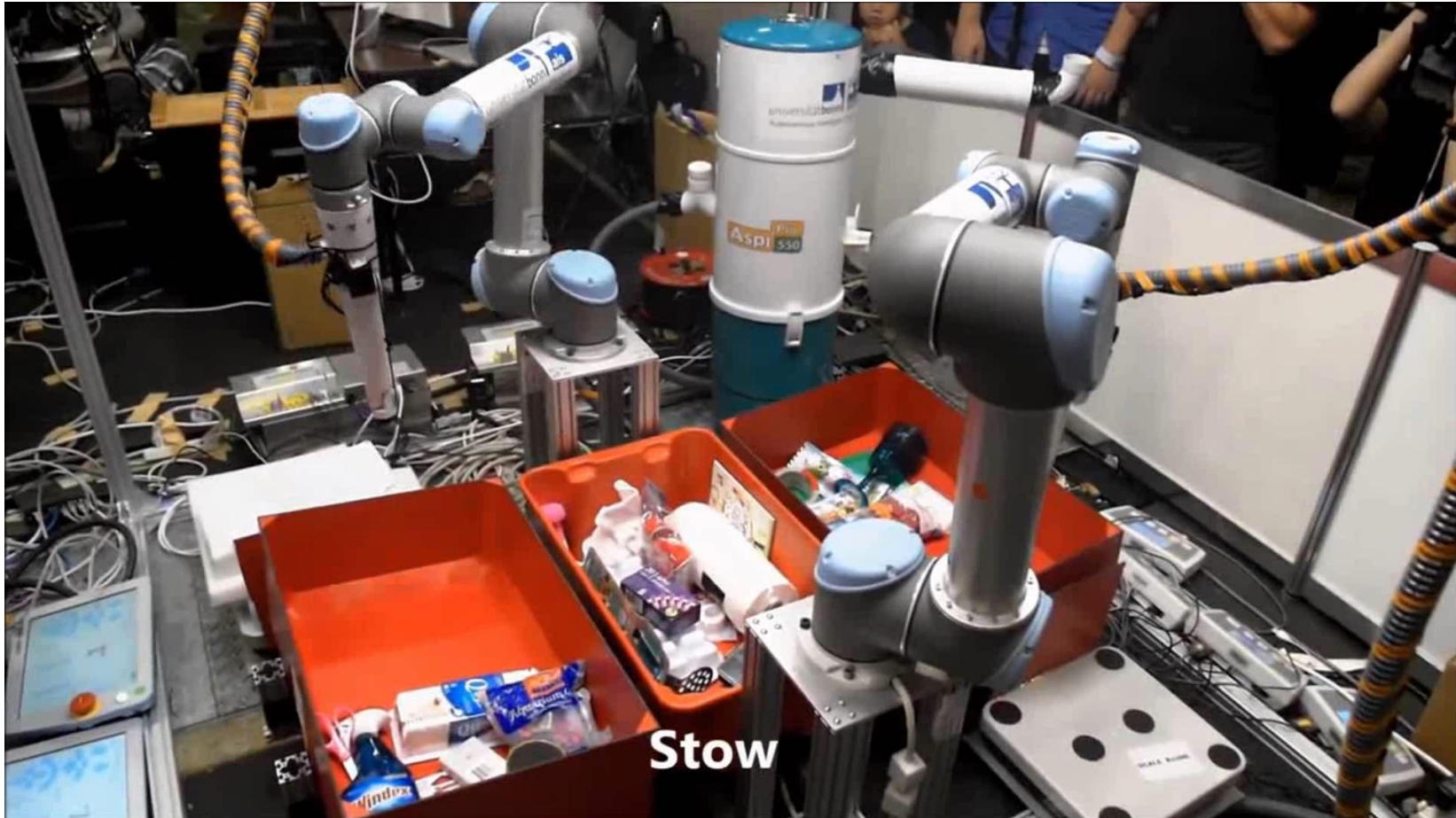
q-tips\_500  
conf: 0.475015

fiskars\_scissors  
conf: 0.831069

ice\_cube\_tray  
conf: 0.976856

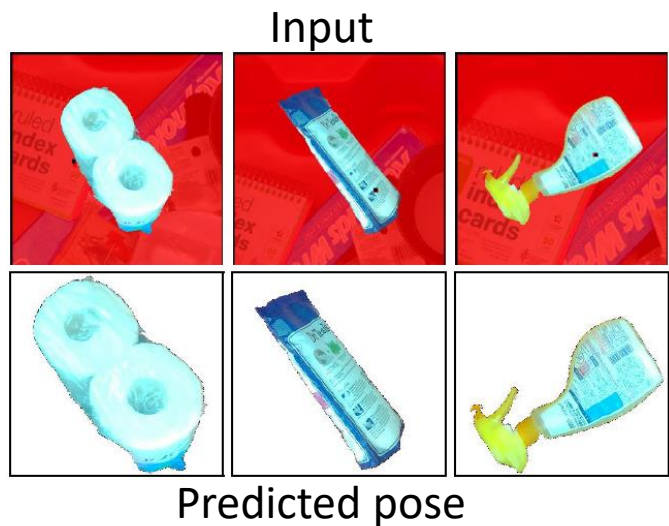
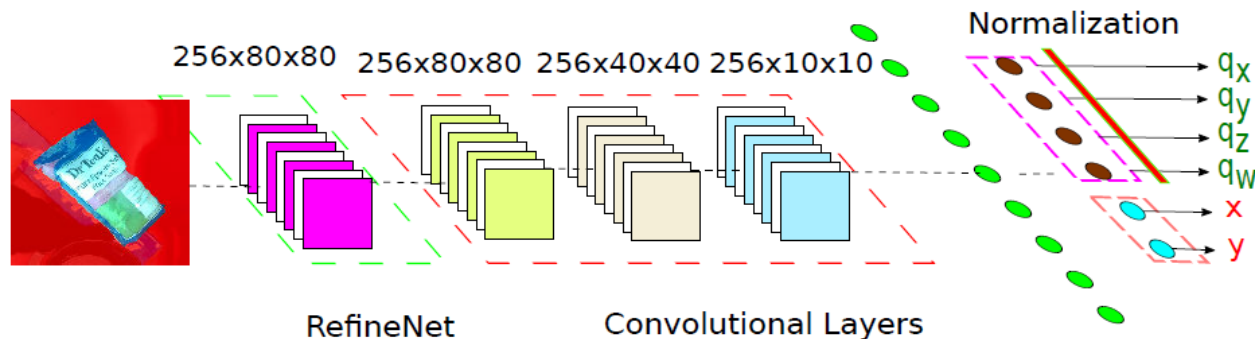


# Amazon Robotics Challenge 2017



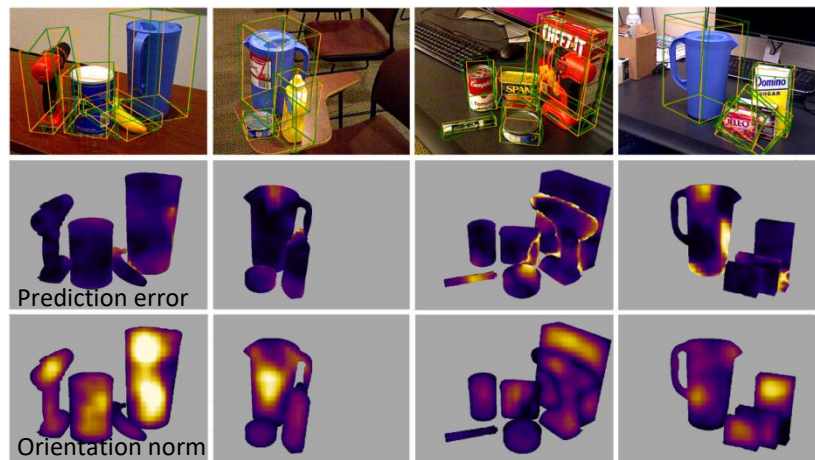
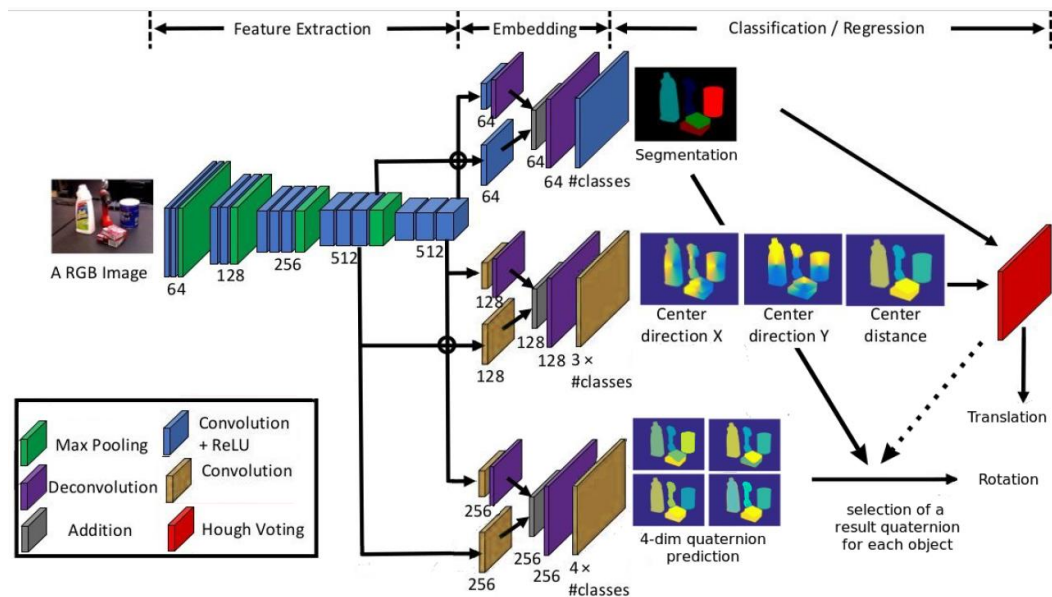
# Object Pose Estimation

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



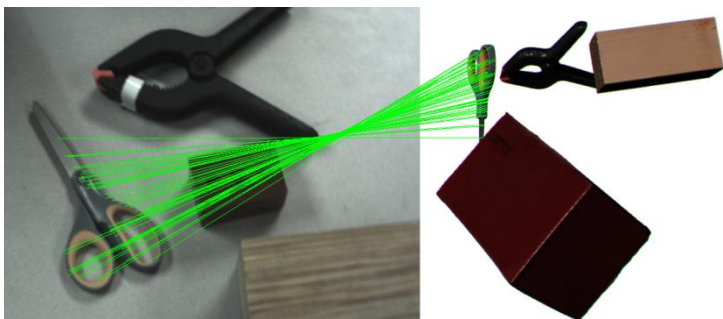
# Dense Convolutional 6D Object Pose Estimation

- Extension of PoseCNN [Xiang et al. RSS 2018]
- Dense prediction of object center and orientation, without cutting out

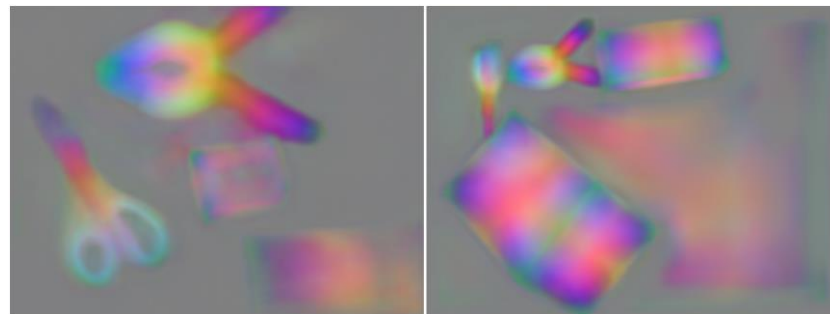


# Self-Supervised Surface Descriptor Learning

- Feature descriptor should be constant under different transformations, viewing angles, and environmental effects such as lighting changes
- Descriptor should be unique to facilitate matching across different frames or representations
- Learn dense features using a contrastive loss



Known correspondences



Learned features

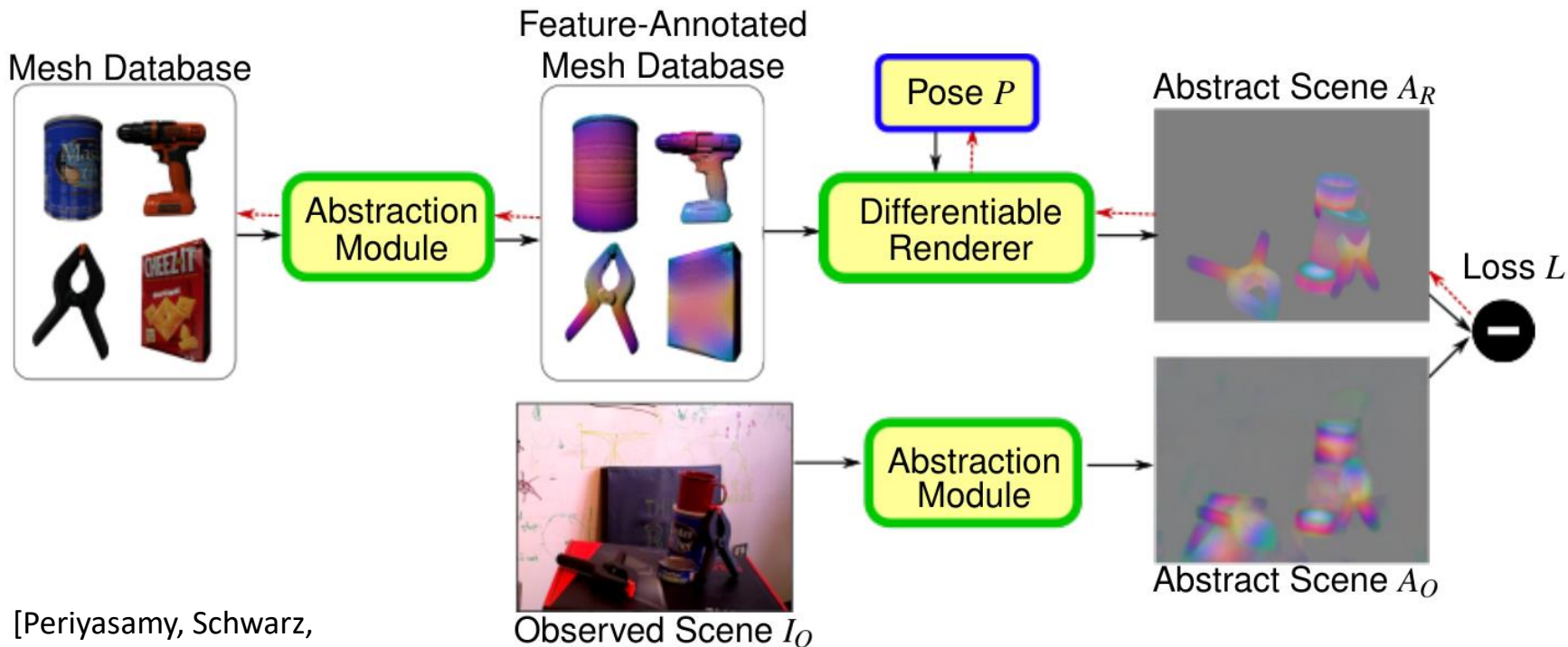
# Descriptors as Texture on Object Surfaces

- Learned feature channels used as textures for 3D object models
- Used for 6D object pose estimation



# Abstract Object Registration

- Compare rendered and actual scene in feature space
- Adapt model pose by gradient descent



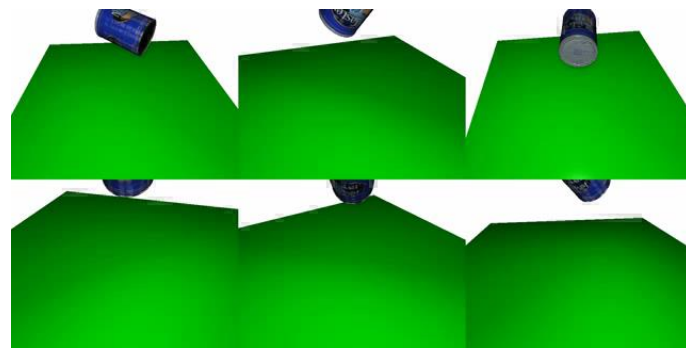
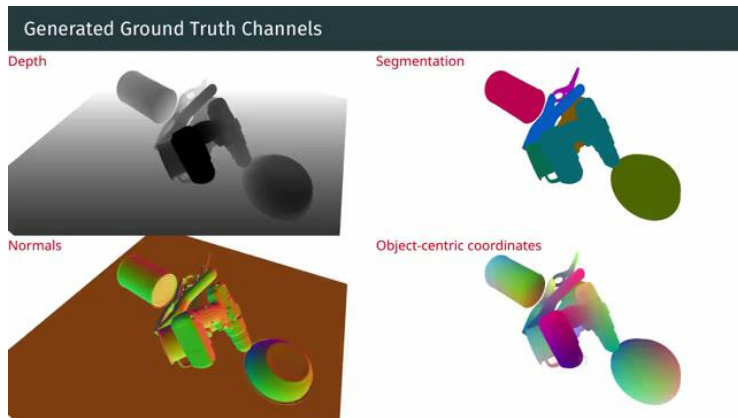
[Periyasamy, Schwarz,  
Behnke Humanoids 2019]

# Registration Examples



# Stilleben: Learning from Synthetic Scenes

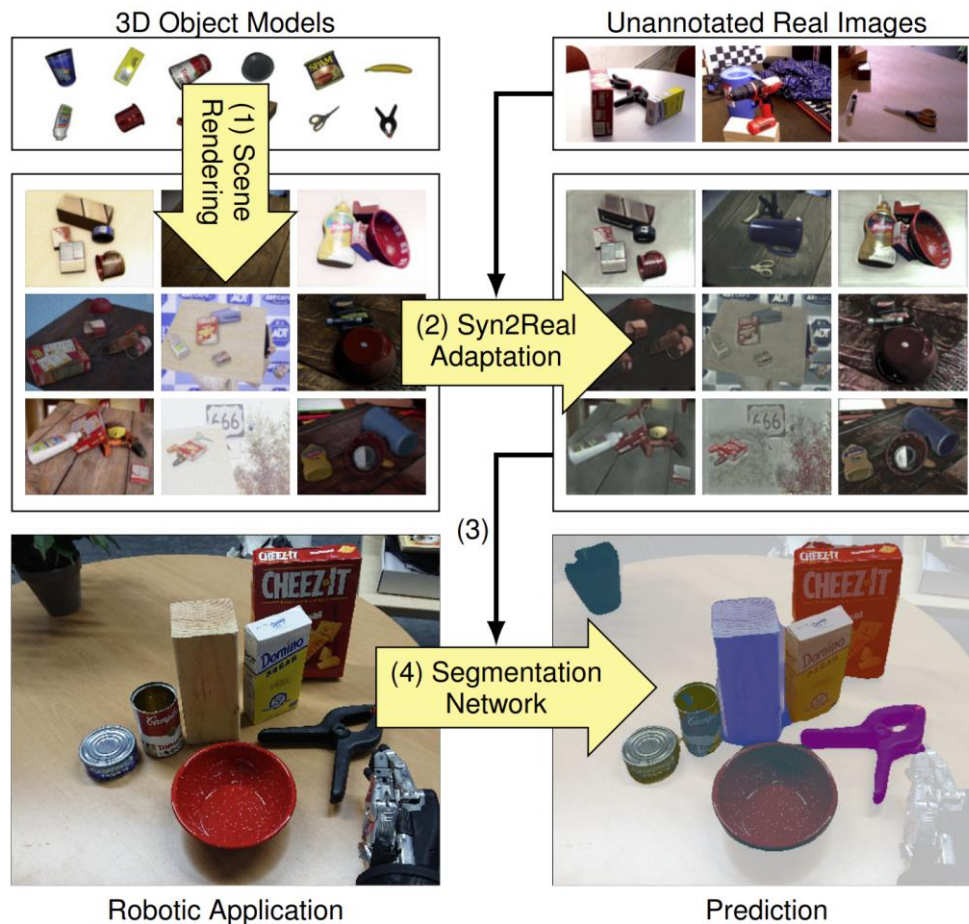
- Cluttered arrangements from 3D meshes
- Photorealistic scenes with randomized material and lighting including ground truth
- For online learning & render-and-compare
- Semantic segmentation on YCB Video Dataset
  - Close to real-data accuracy
  - Improves segmentation of real data





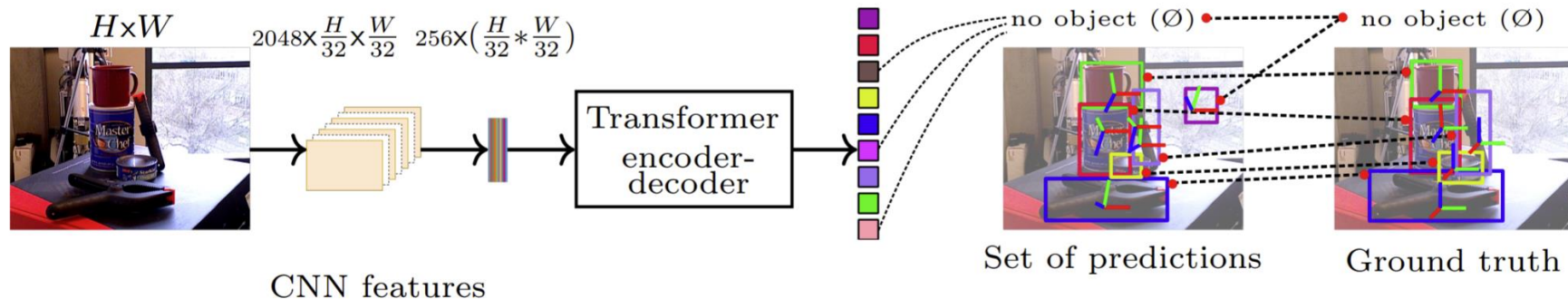
# Synthetic-to-Real Domain Adaptation

- Generate images from 3D object meshes
- Adapt the synthetic images to the real domain using unannotated real images (GAN loss)
- Train downstream task using adapted images
- Semantic segmentation results almost as good as trained with real images
- Improved results in combination with real annotations

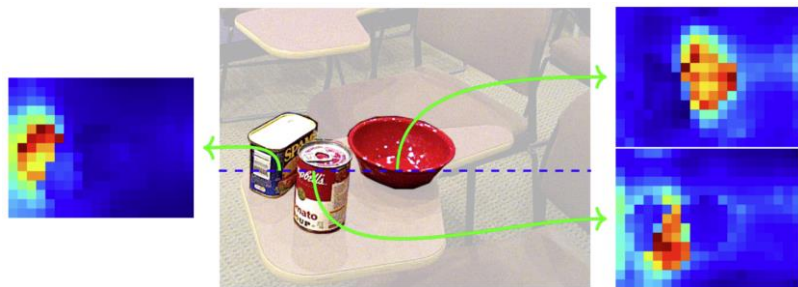


# T6D-Direct: Transformers for Multi-Object 6D Pose Direct Regression

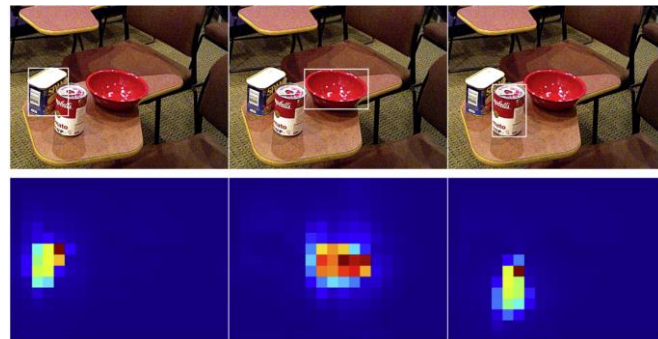
- Extends DETR: End-to-end object detection with transformers [Carion et al. ECCV 2020]
- End-to-end differentiable pipeline for 6D object pose estimation



Encoder self-attention

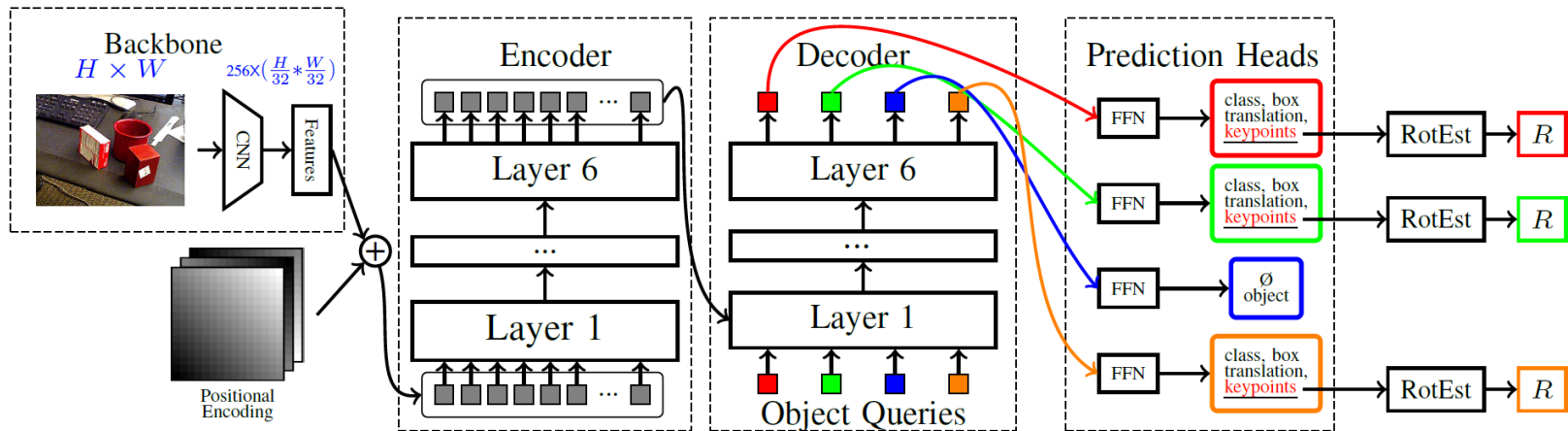


Object detections and decoder attention



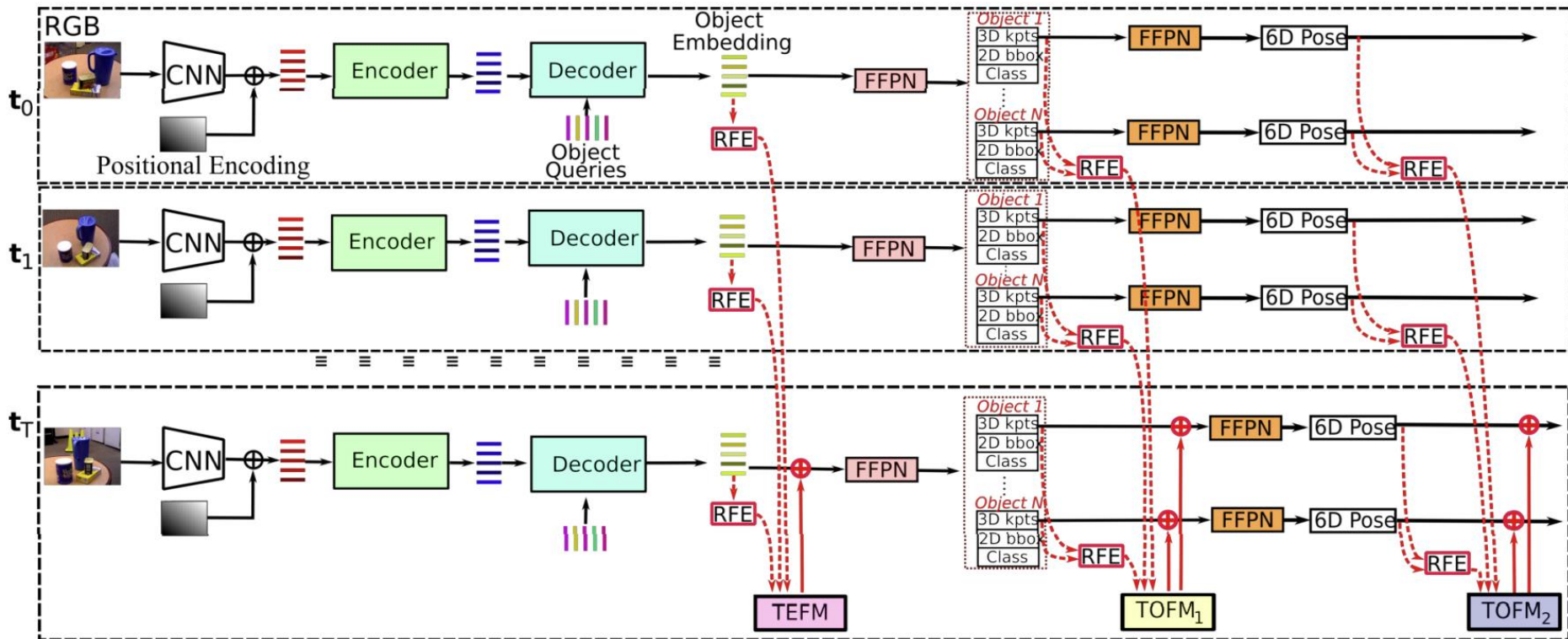
[Amini et al. GCPR 2021]

# Multi-Object 6D Pose Estimation using Keypoint Regression



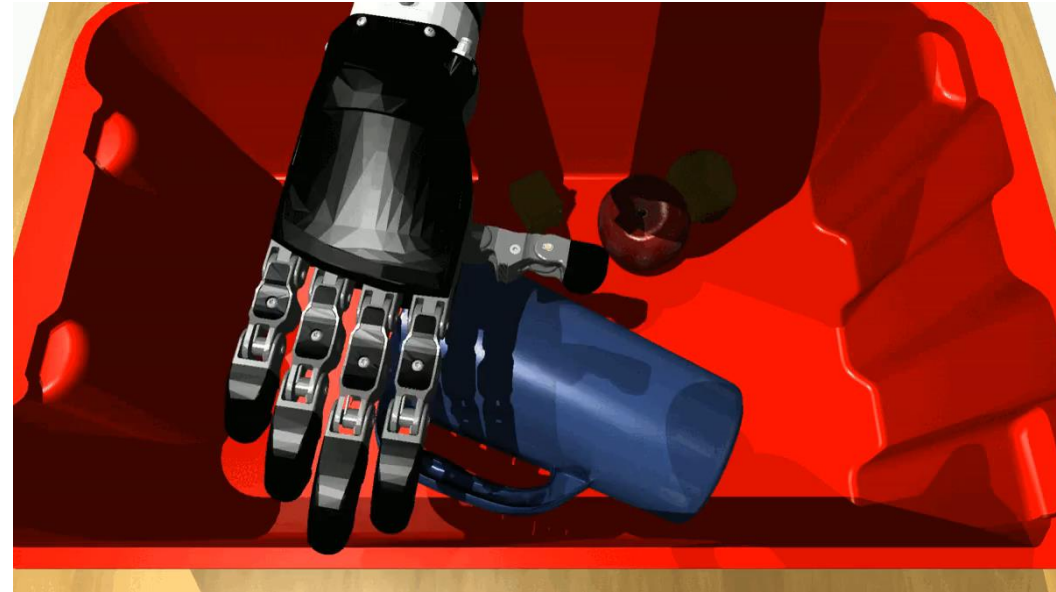
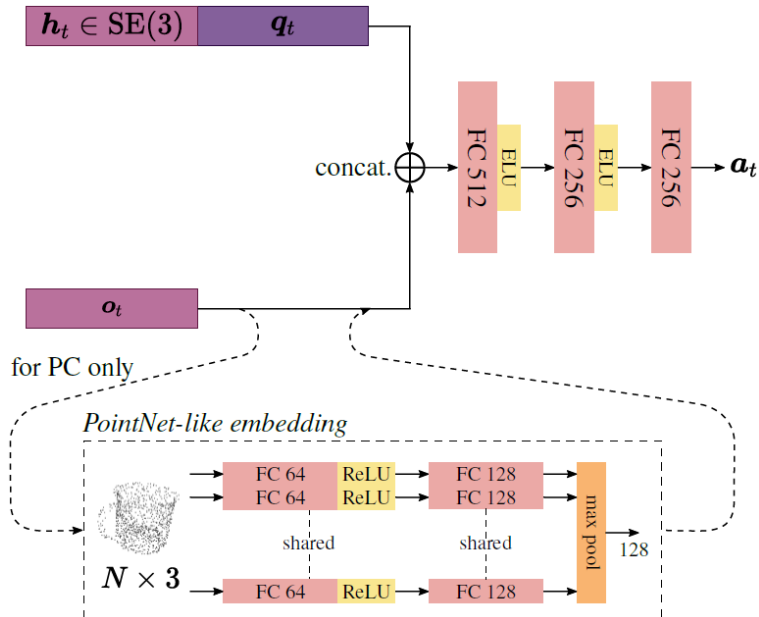
# MOTPose: Attention-based Temporal Fusion for Multi-object 6D Pose Estimation

- Transformer-based temporal embedding and object fusion modules



# Learning Interactive Grasping

- Deep RL-based interactive policy
- Input: object parameters or point cloud + hand pose
- Output: increments of hand DoF:

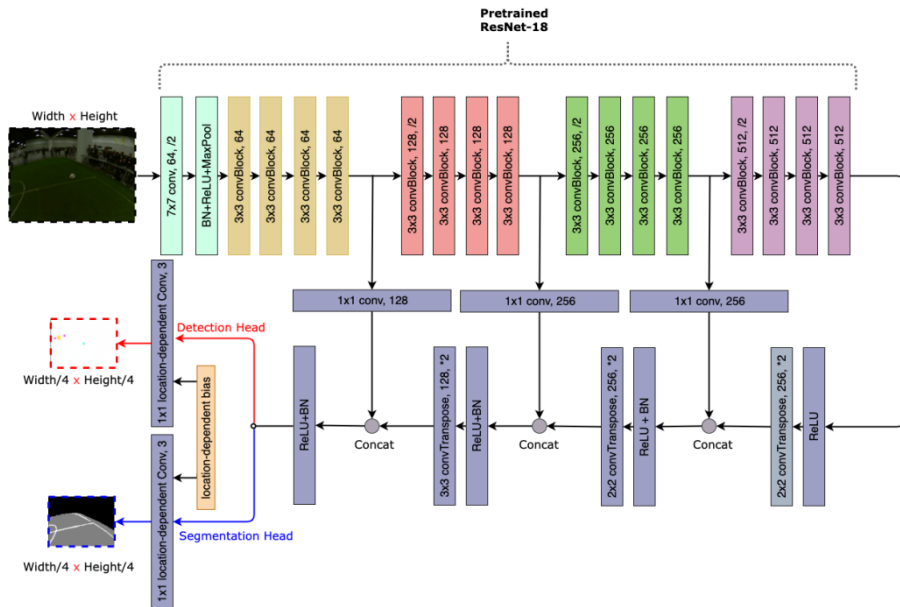


# Humanoid AdultSize Soccer: RoboCup 2022 in Bangkok



# Transfer Learning for Visual Perception

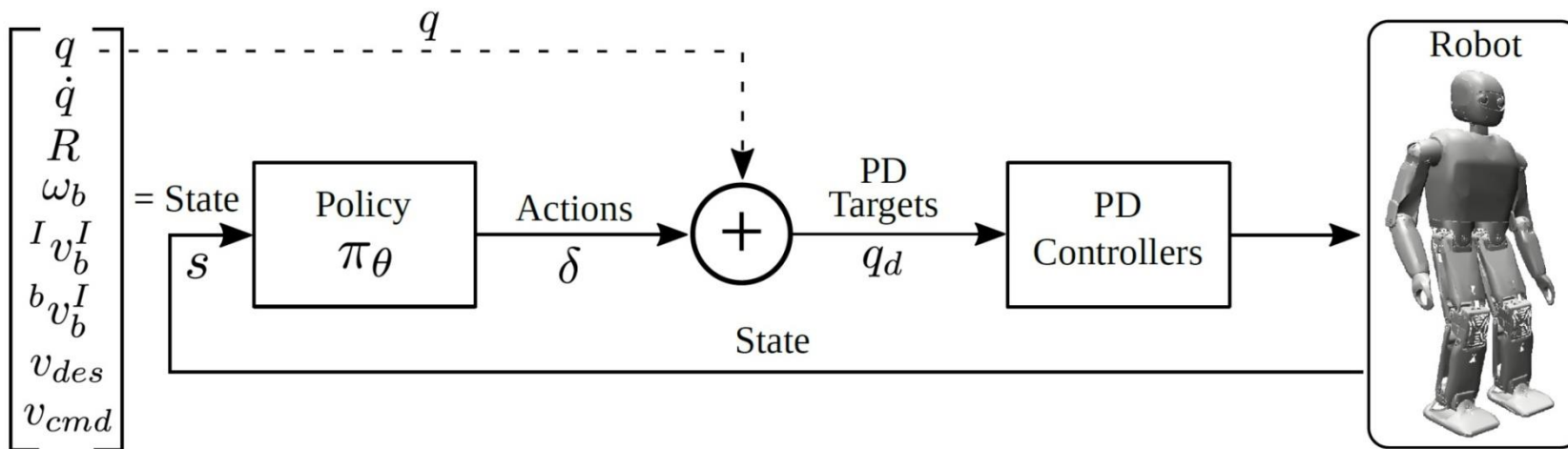
- Encoder-decoder network
- Two outputs
  - Object detection
  - Semantic segmentation
- Location-dependent bias



- Detects objects that are hard to recognize for humans
- Robust to lighting changes

# Learning Omnidirectional Gait from Scratch

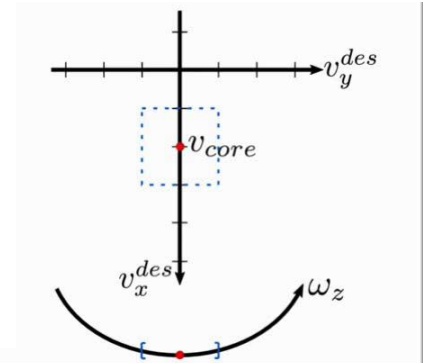
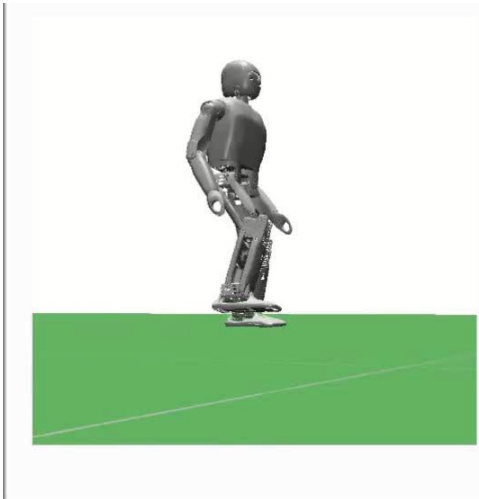
- State includes joint positions and velocities, robot orientation, robot speed
- Actions are increments of joint positions
- Simple reward structure
  - Velocity tracking
  - Pose regularization
  - Not falling





# Learning Curriculum

- Start with small velocities
- Increase range of sampled velocities



# Learned Omnidirectional Gait

- Target velocity can be changed continuously

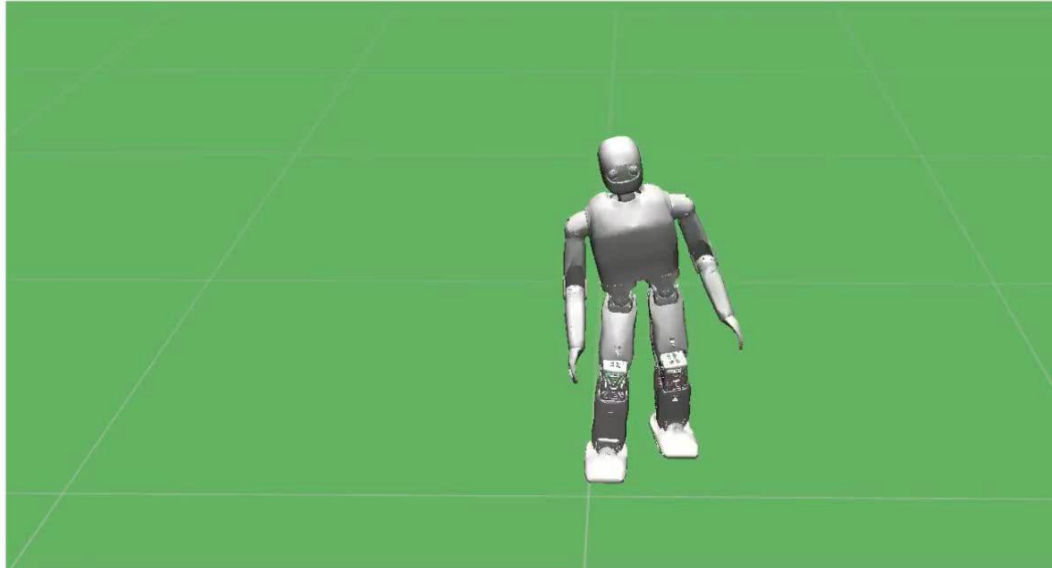
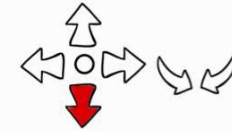
Our locomotion controller is able to:

**Walk Forward**

$$v_x = 0.6 \text{ m/s}$$

$$v_y = 0.0 \text{ m/s}$$

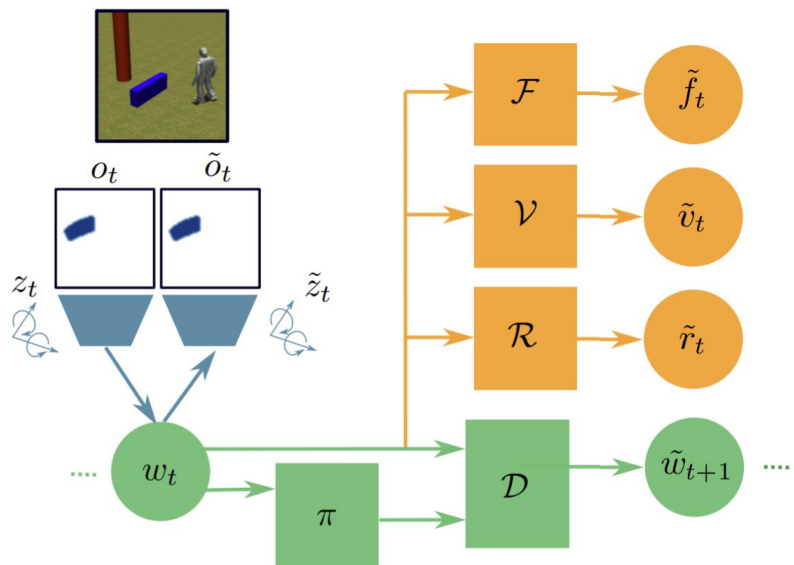
$$\omega_z = 0.0 \text{ rad/s}$$



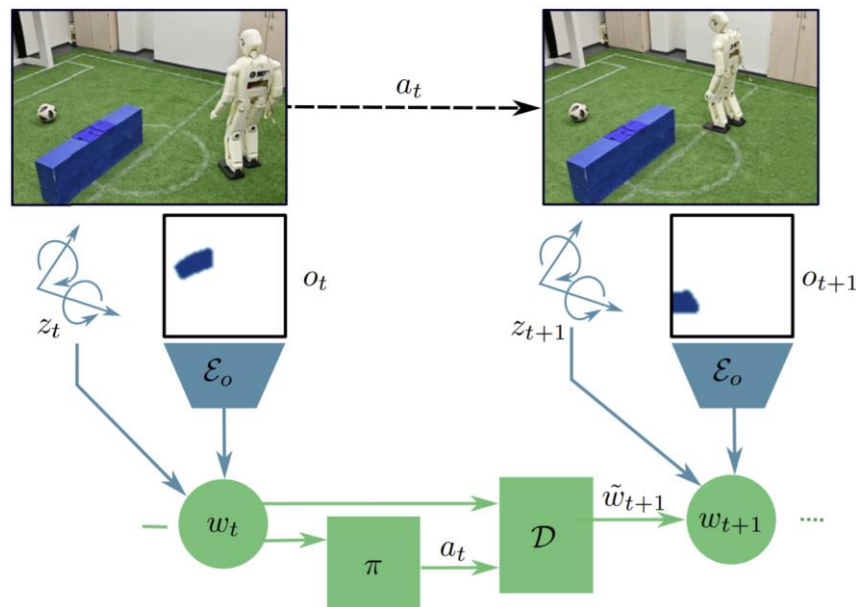
# Learning Mapless Humanoid Navigation

- Visual (RGB images) and nonvisual observations to learn a control policy and an environment dynamics model
- Anticipate terminal states of success and failure

Training



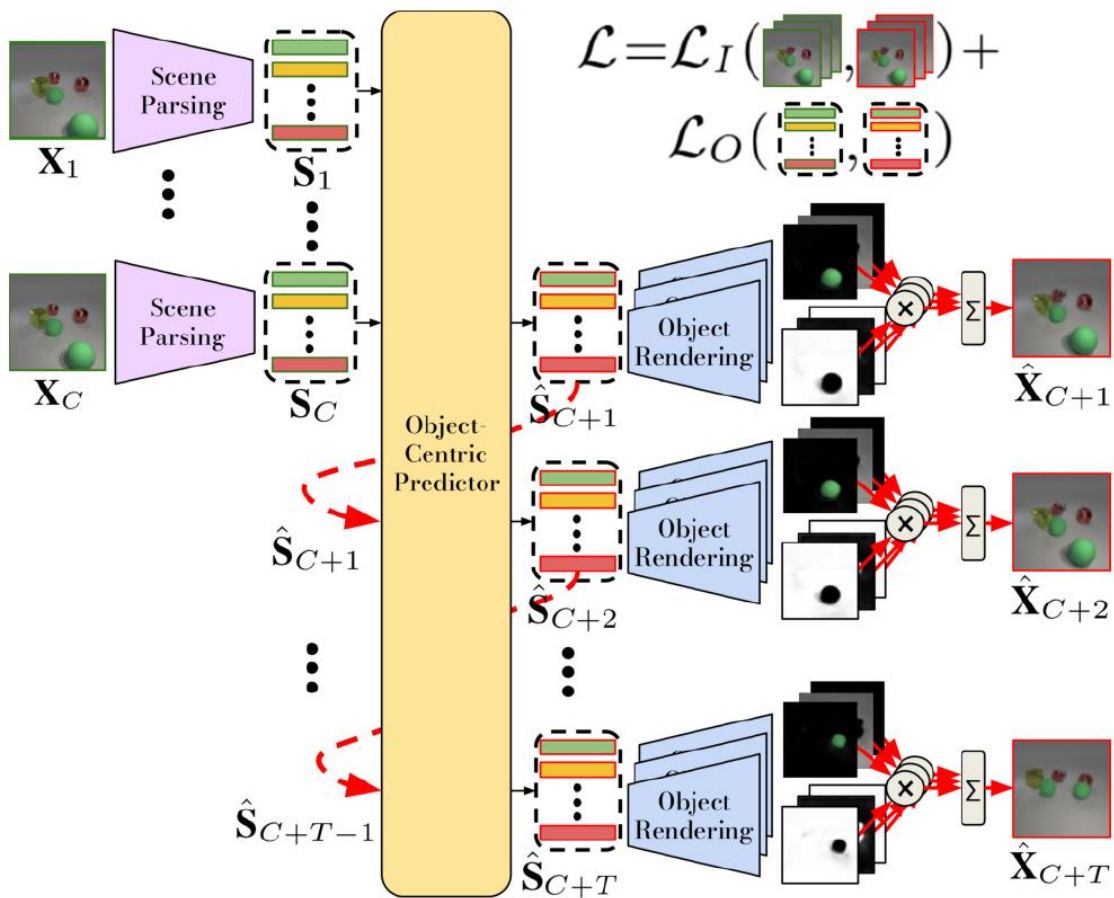
Inference



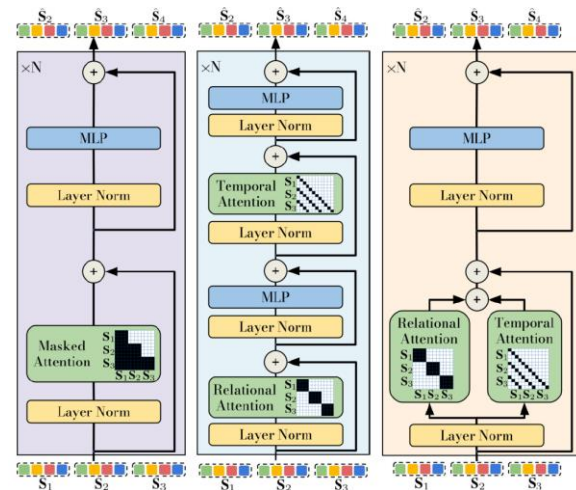
# Learning Mapless Humanoid Navigation



# Object-centric Video Prediction Decoupling Dynamics and Interaction



- Scene parsing into object slots
- Video synthesis from objects and masks
- Predictor decouples temporal and relational attention



# Object-centric Video Prediction Decoupling Dynamics and Interaction

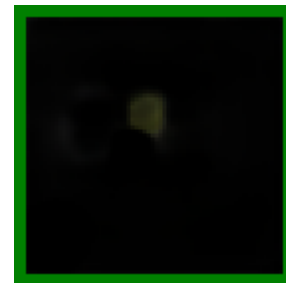
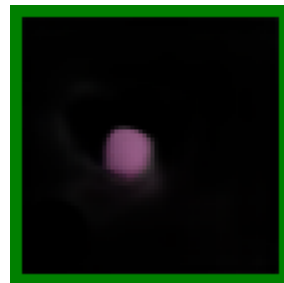
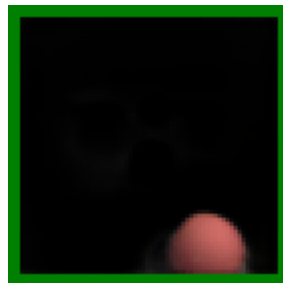
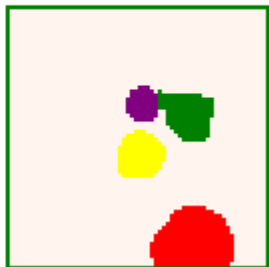
**GT**



**Pred. RGB**



**Pred. Masks**



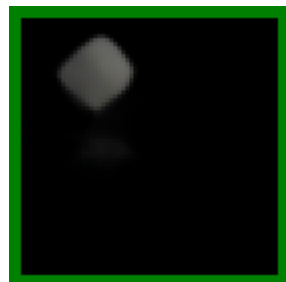
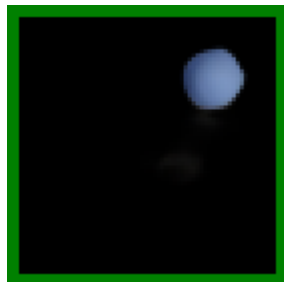
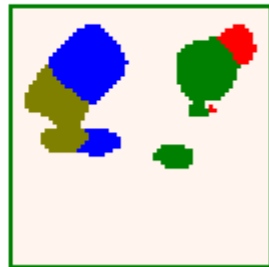
**GT**



**Pred. RGB**

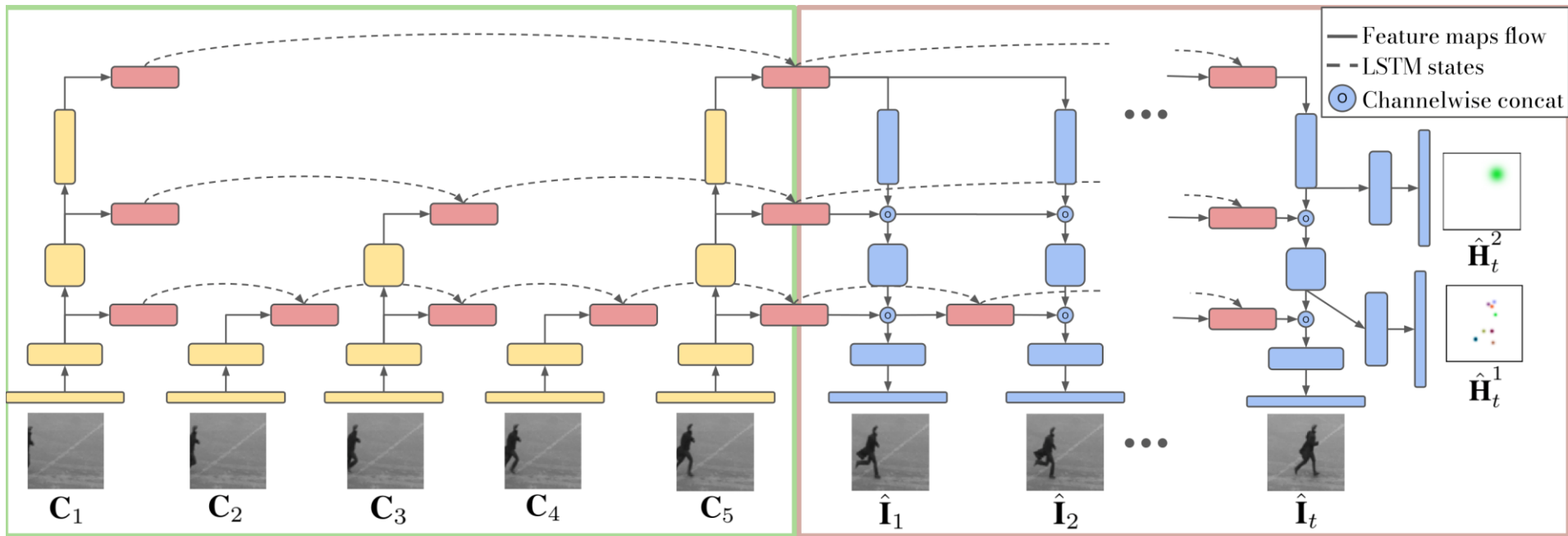


**Pred. Masks**



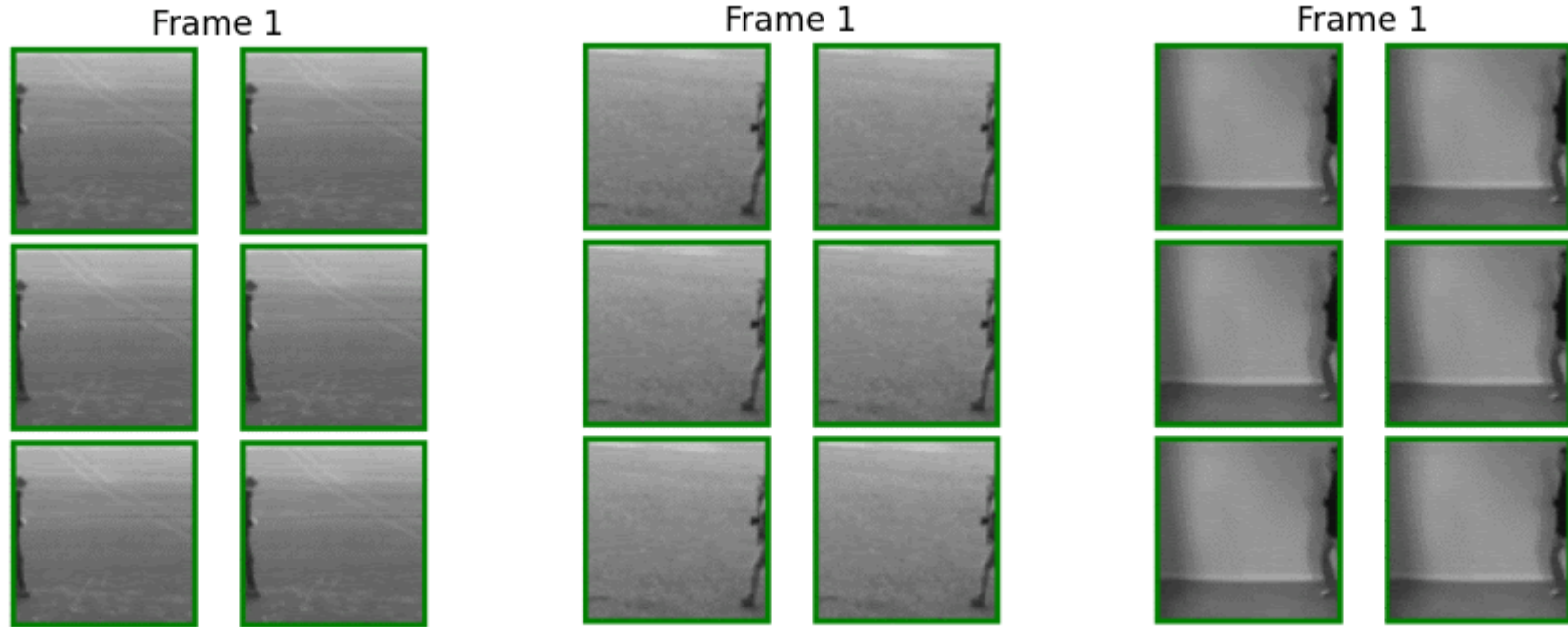
# MSPred: Video Prediction at Multiple Spatio-Temporal Scales

- Coarser, more abstract predictions for longer time horizons in higher layers
- Predict image itself, human pose joint keypoints, and human body position



# MSPred: Video Prediction at Multiple Spatio-Temporal Scales

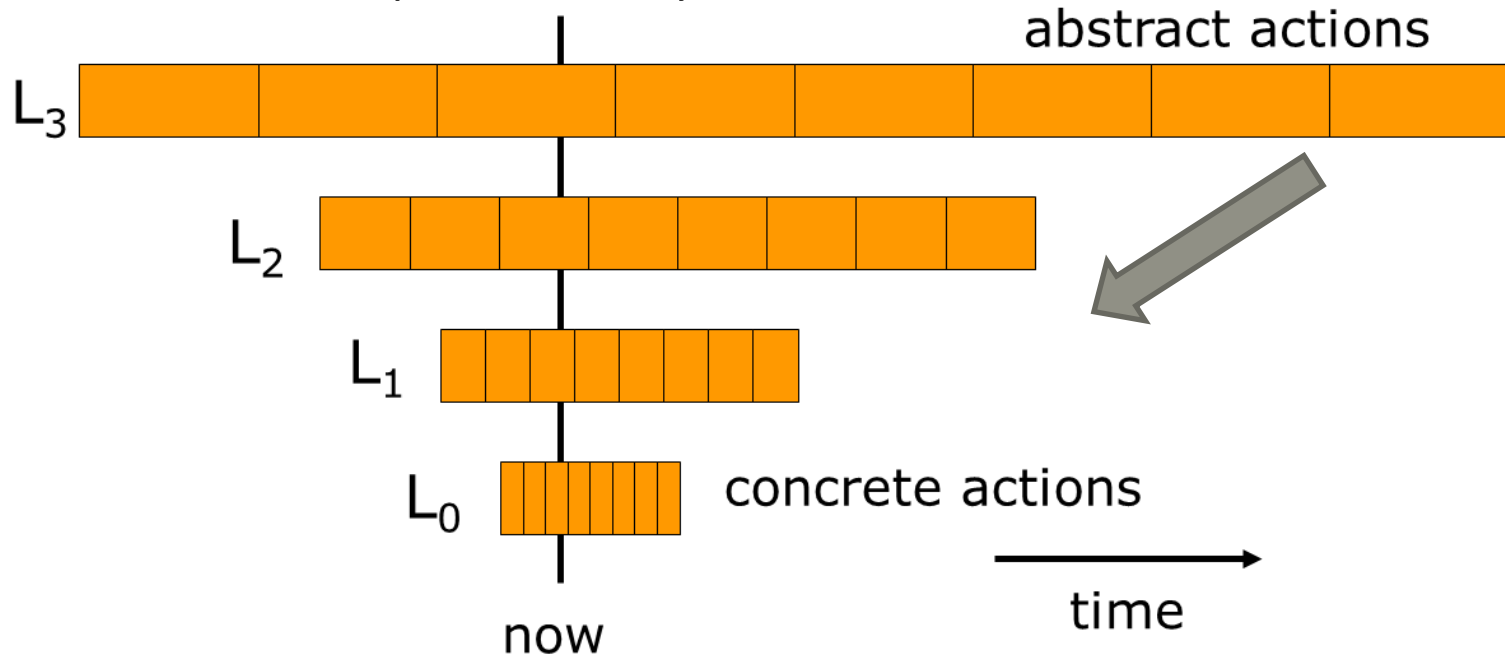
- Coarser, more abstract predictions for longer time horizons in higher layers
- Predict image itself, human pose joint keypoints, and human body position



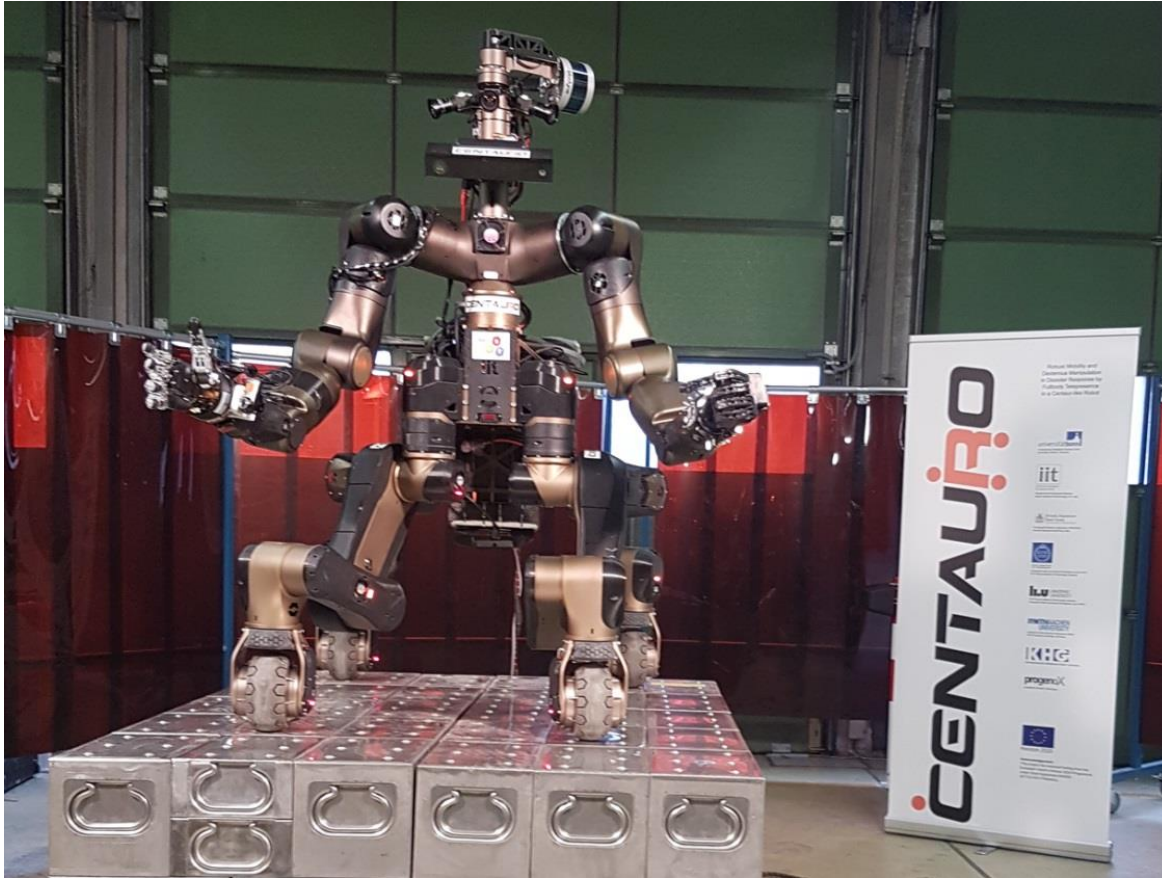


# Hierarchical Planning in the Now

- Use predicted state on different layers of abstraction for planning
- Coarse-to-fine planning makes actions more concrete as they come closer to execution
- Plan consists of few steps on each layer



# Centauro Robot



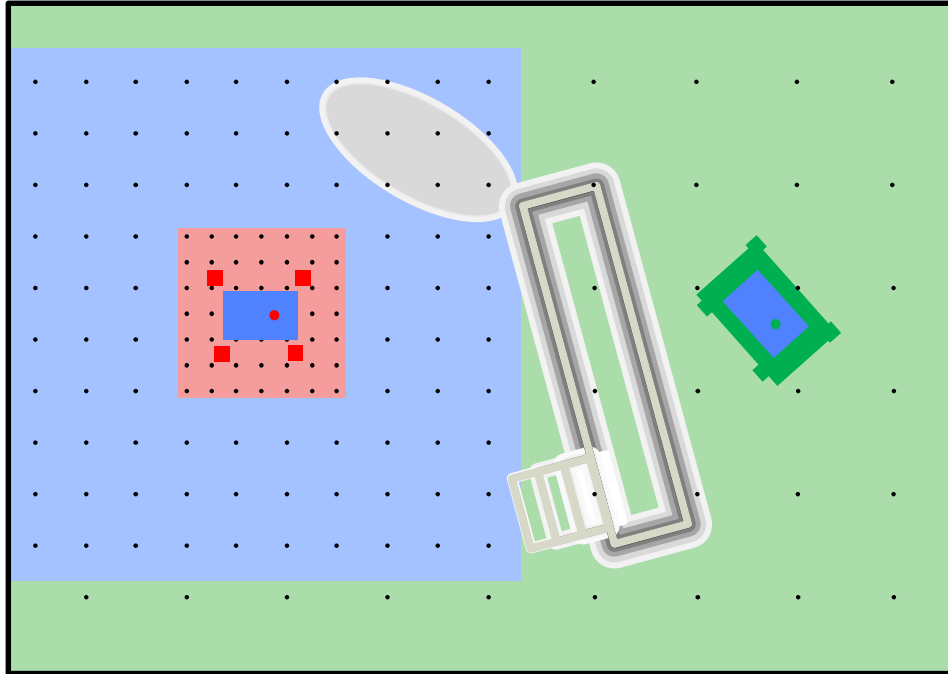
# CENTAURO

- Serial elastic actuators
- 42 main DoFs
- Schunk hand
- 3D laser
- RGB-D camera
- Color cameras
- Two GPU PCs

[Tsagarakis et al., IIT 2017]

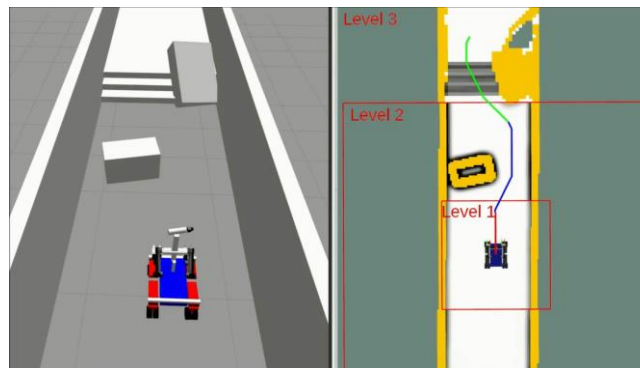
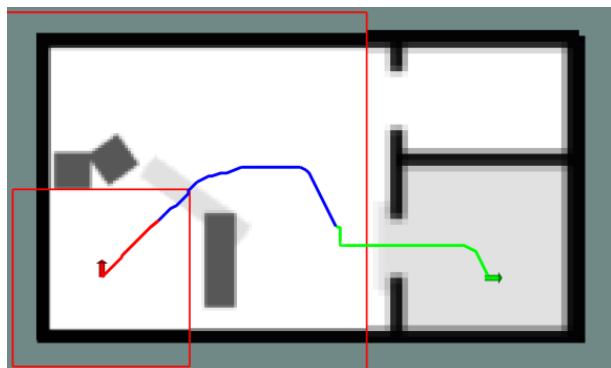
# Hybrid Driving-Stepping Locomotion Planning: Abstraction

- Planning in the here and now
- Far-away details are abstracted away



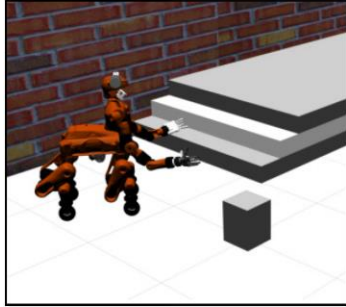
# Hybrid Driving-Stepping Locomotion Planning: Abstraction

Level	Map Resolution	Map Features	Robot Representation	Action Semantics
1	<ul style="list-style-type: none"> <li>• 2.5 cm</li> <li>• 64 orient.</li> </ul>	<ul style="list-style-type: none"> <li>• Height</li> </ul>		<ul style="list-style-type: none"> <li>• Individual Foot Actions</li> </ul>
2	<ul style="list-style-type: none"> <li>• 5.0 cm</li> <li>• 32 orient.</li> </ul>	<ul style="list-style-type: none"> <li>• Height</li> <li>• Height Difference</li> </ul>		<ul style="list-style-type: none"> <li>• Foot Pair Actions</li> </ul>
3	<ul style="list-style-type: none"> <li>• 10 cm</li> <li>• 16 orient.</li> </ul>	<ul style="list-style-type: none"> <li>• Height</li> <li>• Height Difference</li> <li>• Terrain Class</li> </ul>		<ul style="list-style-type: none"> <li>• Whole Robot Actions</li> </ul>



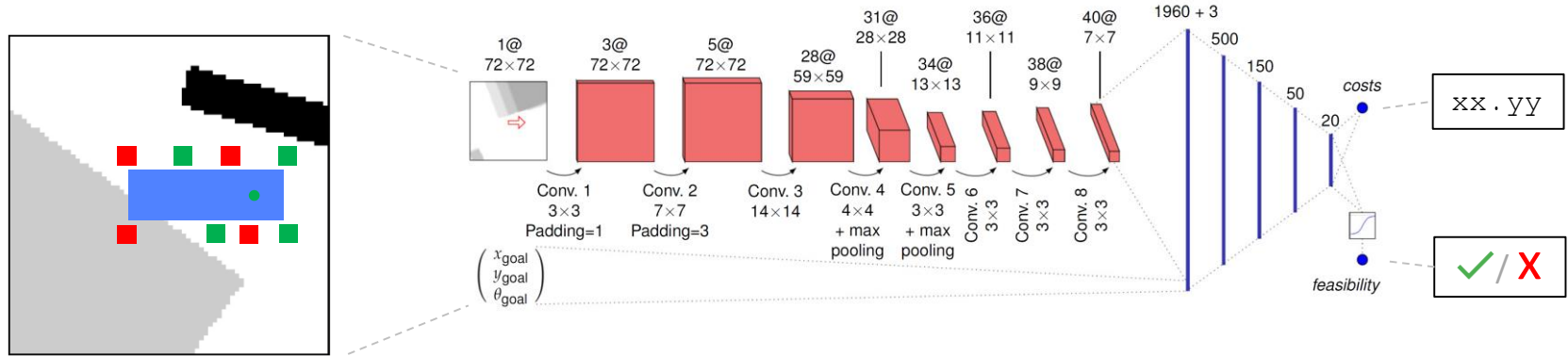
# Learning Cost Functions of Abstract Representations

Planning problem



# Abstraction CNN

- Predict feasibility and costs of local detailed planning

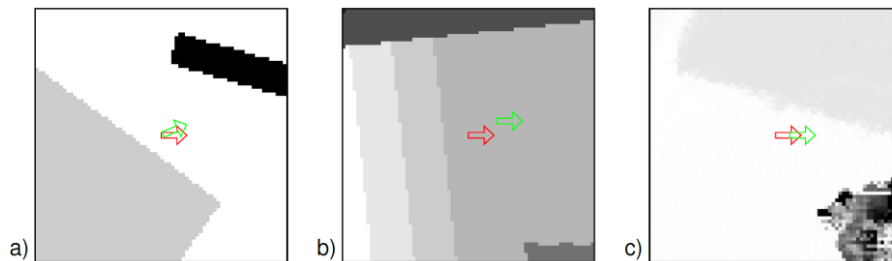


## Training data

- generated with random obstacles, walls, staircases
- *costs* and *feasibility* from detailed A\*-planner
- ~250.000 tasks

# Learned Cost Function: Abstraction Quality

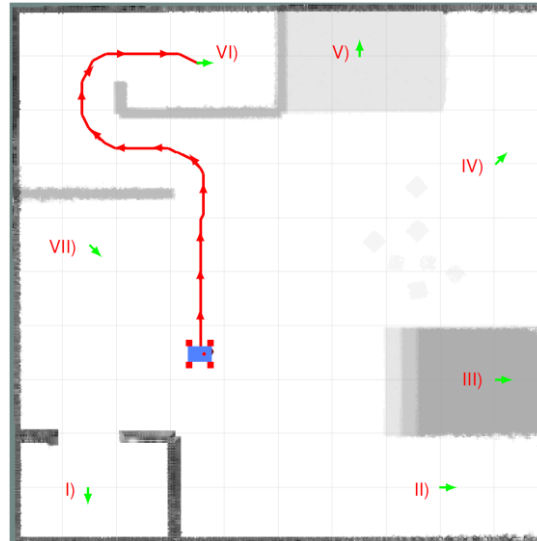
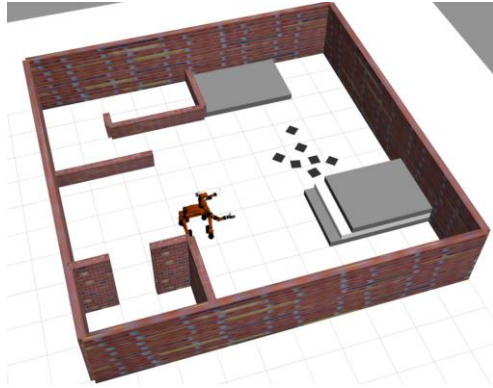
- CNN predicts feasibility and costs better than manually tuned geometric heuristics



	<i>random</i>	<i>simulated</i>	<i>real</i>
<i>feasibility correct, man.tuned</i>	79.27%	65.35%	69.77%
$\text{Error}(\mathcal{C}_{a,\text{man.tuned}})$	0.057	0.021	0.103
<i>feasibility correct, CNN</i>	95.04%	96.69%	92.62%
$\text{Error}(\mathcal{C}_{a,\text{CNN}})$	0.027	0.013	0.081

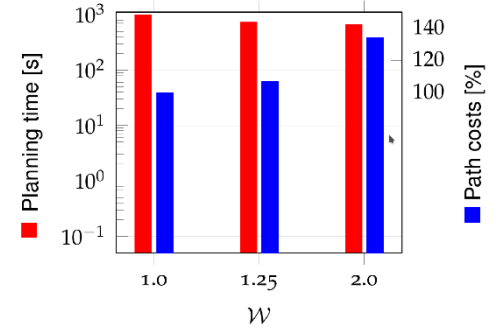
# Experiments – Planning Performance

- Learned heuristics accelerates planning, without increasing path costs much

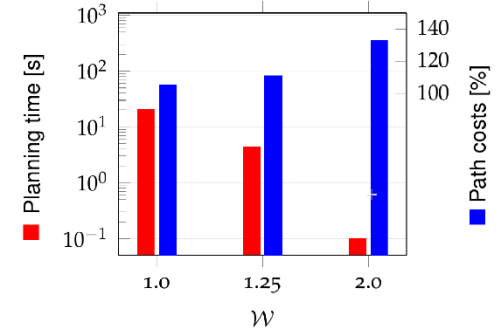


Heuristic preprocessing: 239 sec

### Geometric heuristic



### Abstract representation heuristic





# CENTAURO Evaluation @ KHG: Locomotion Tasks



# Transfer of Manipulation Skills

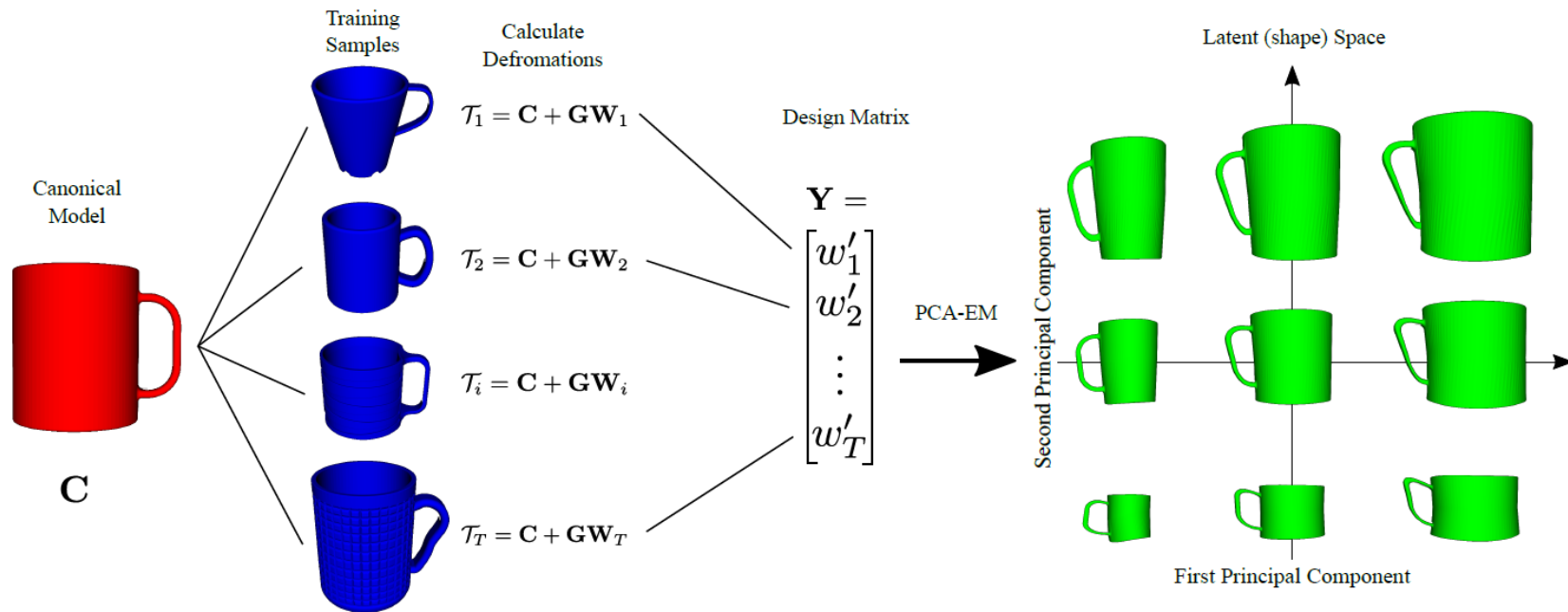


Knowledge  
Transfer

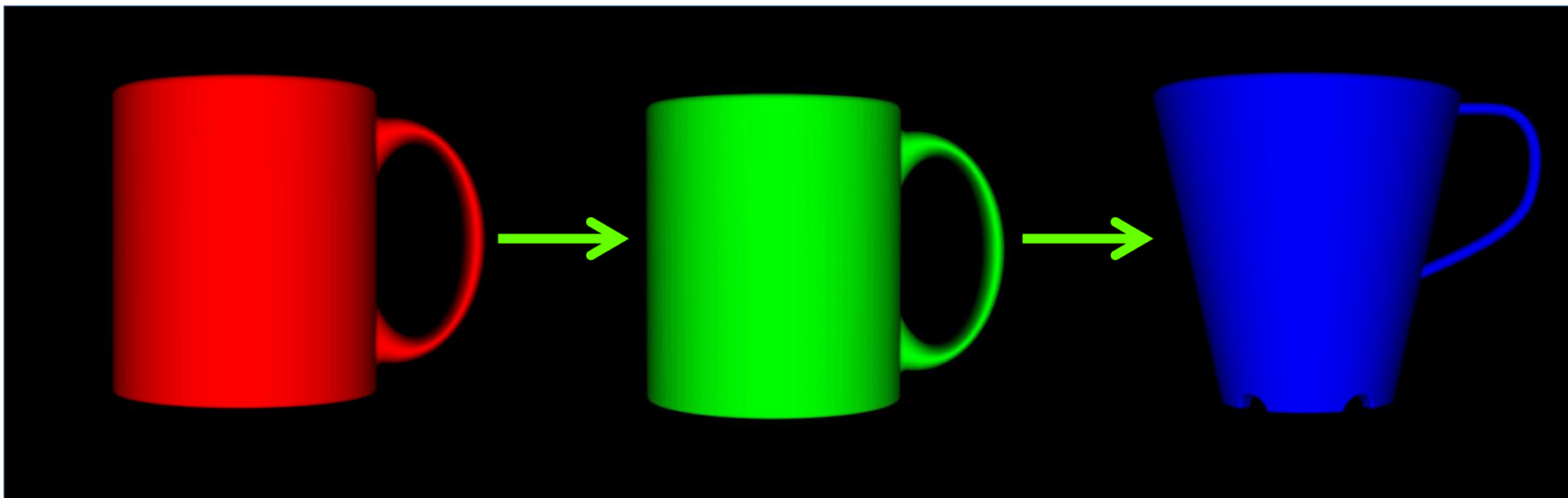


# Learning a Latent Shape Space

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

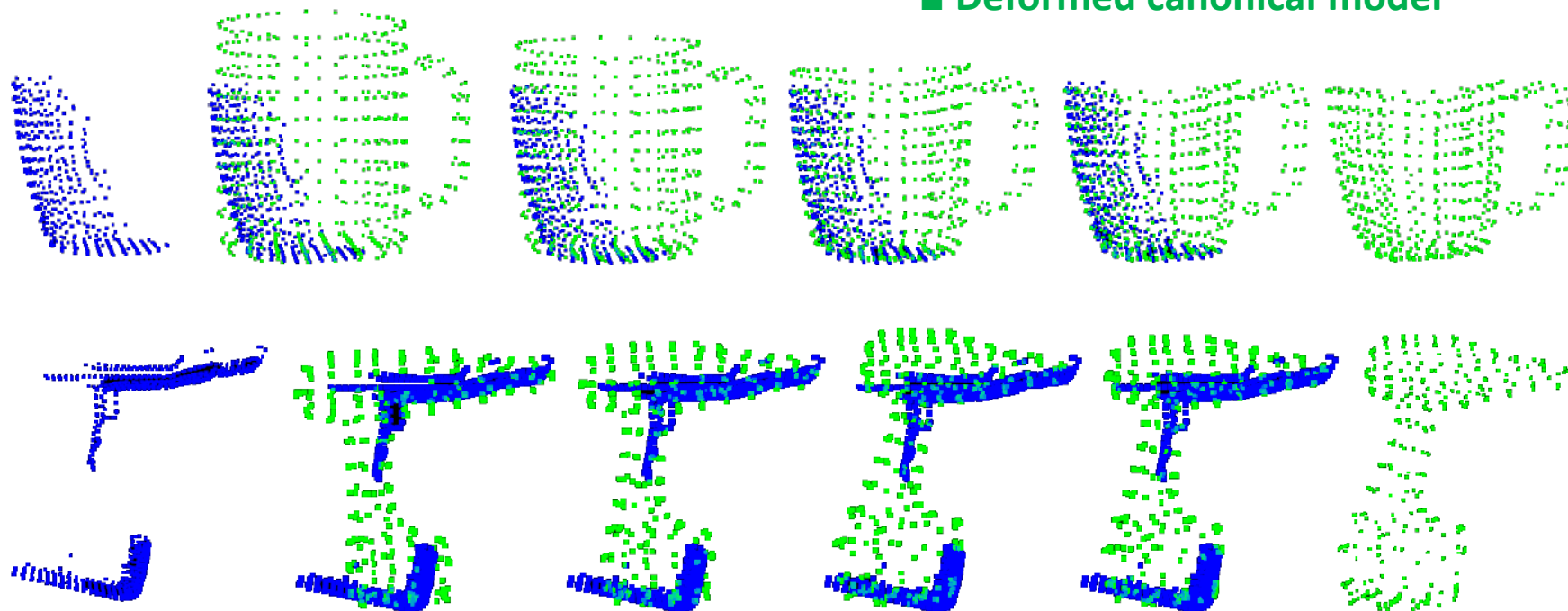


# Interpolation in Shape Space



# Shape-aware Non-rigid Registration

- Partial view of novel instance
- Deformed canonical model

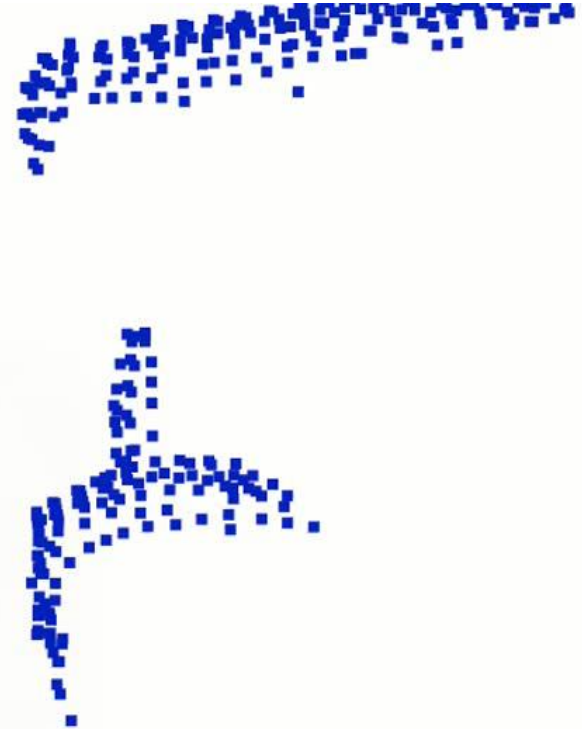


# Shape-aware Registration for Grasp Transfer

■ Full point cloud



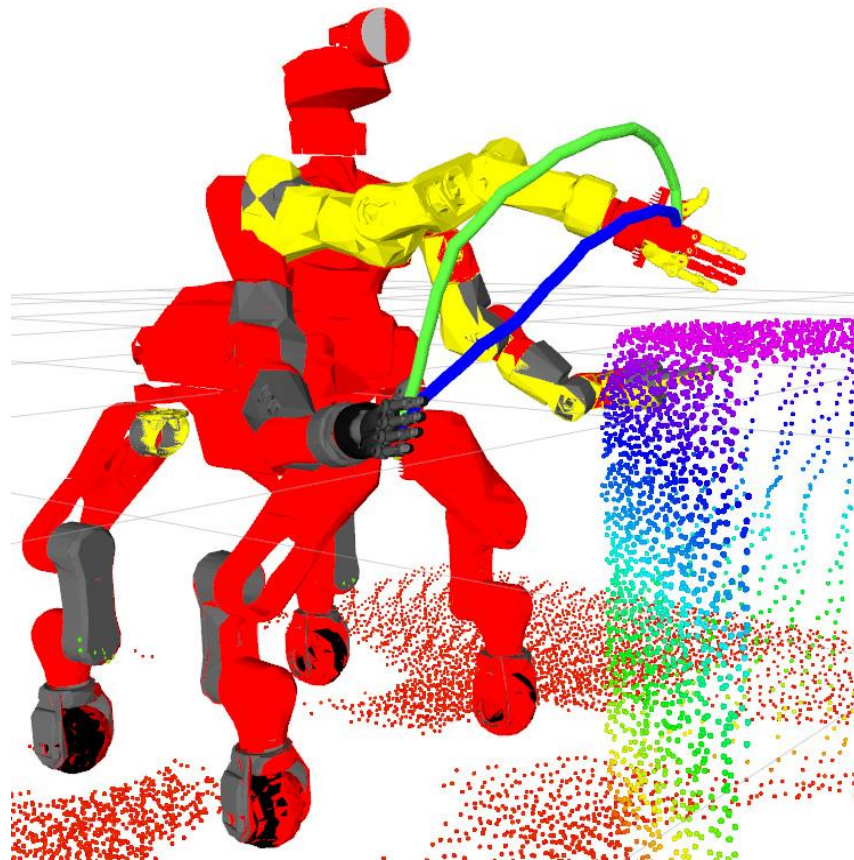
■ Partial view



# Collision-aware Motion Generation

Constrained Trajectory Optimization:

- Collision avoidance
- Joint limits
- Time minimization
- Torque optimization

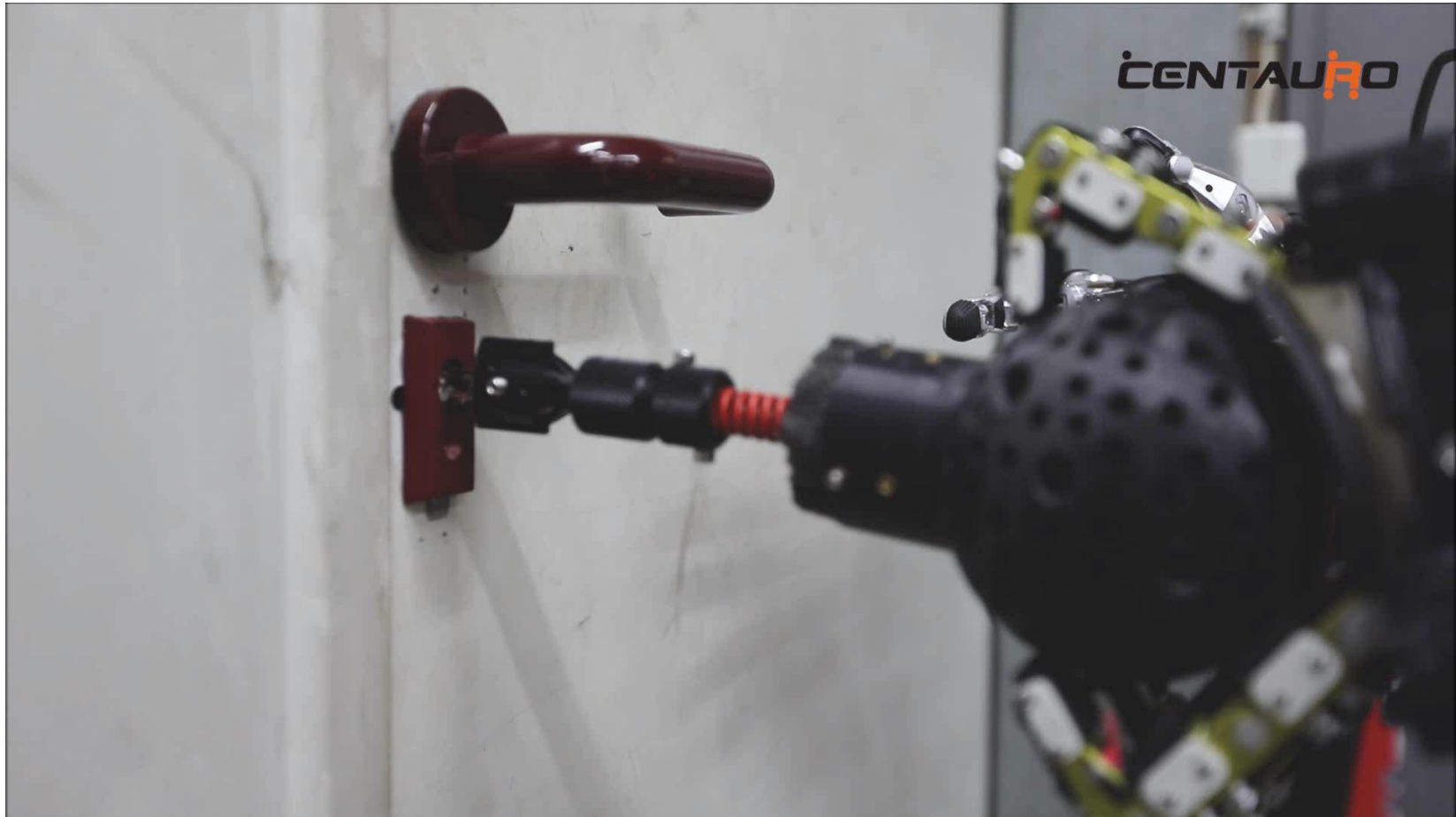


# Grasping an Unknown Power Drill and Fastening Screws



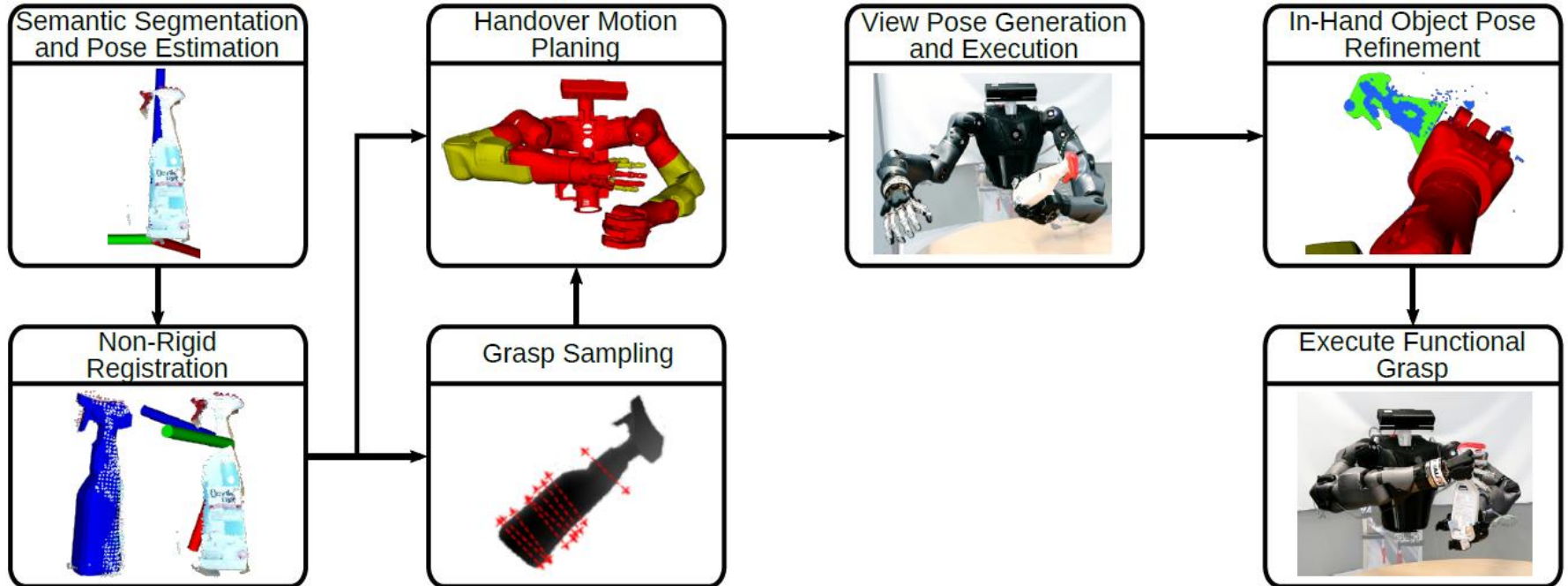


# CENTAURO: Complex Manipulation Tasks

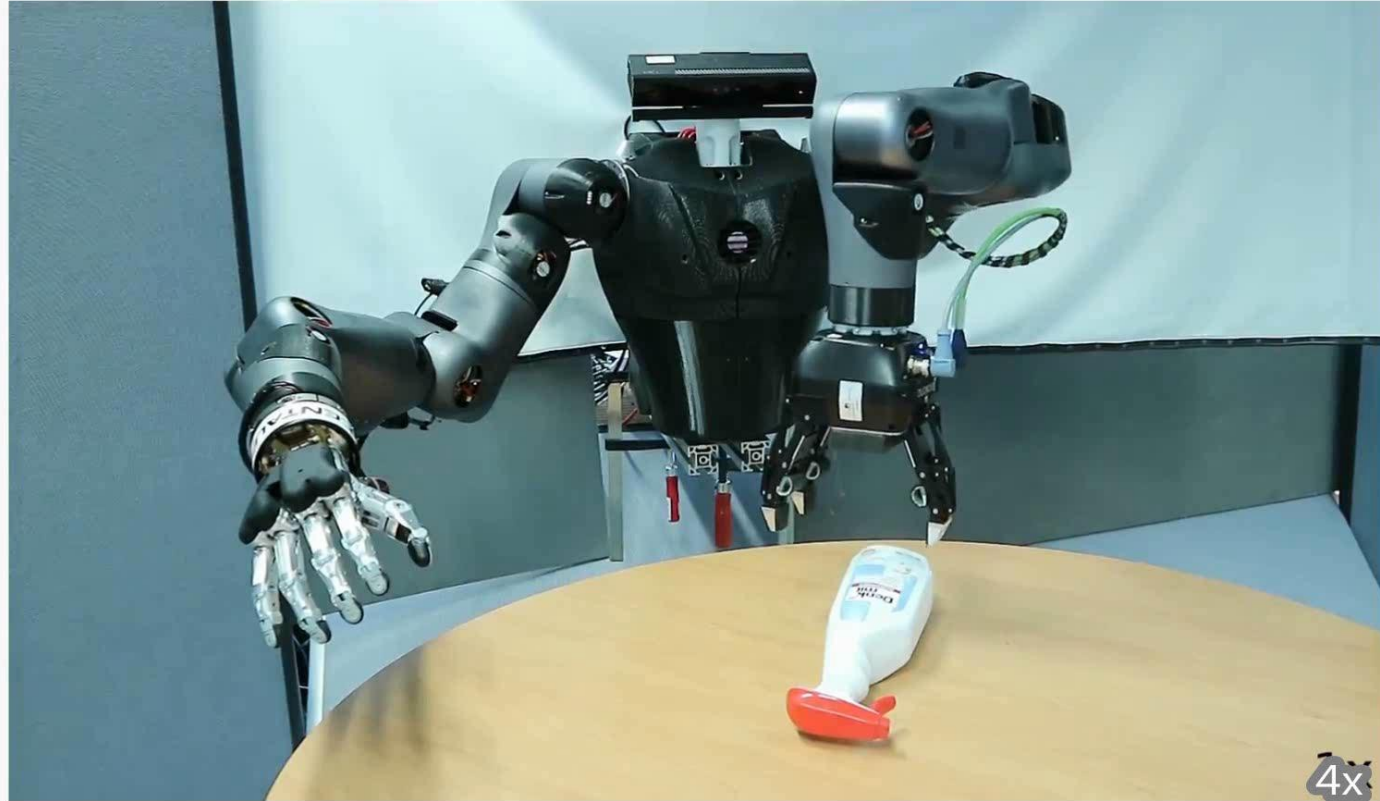


# Regrasping for Functional Grasp

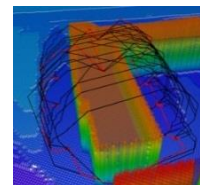
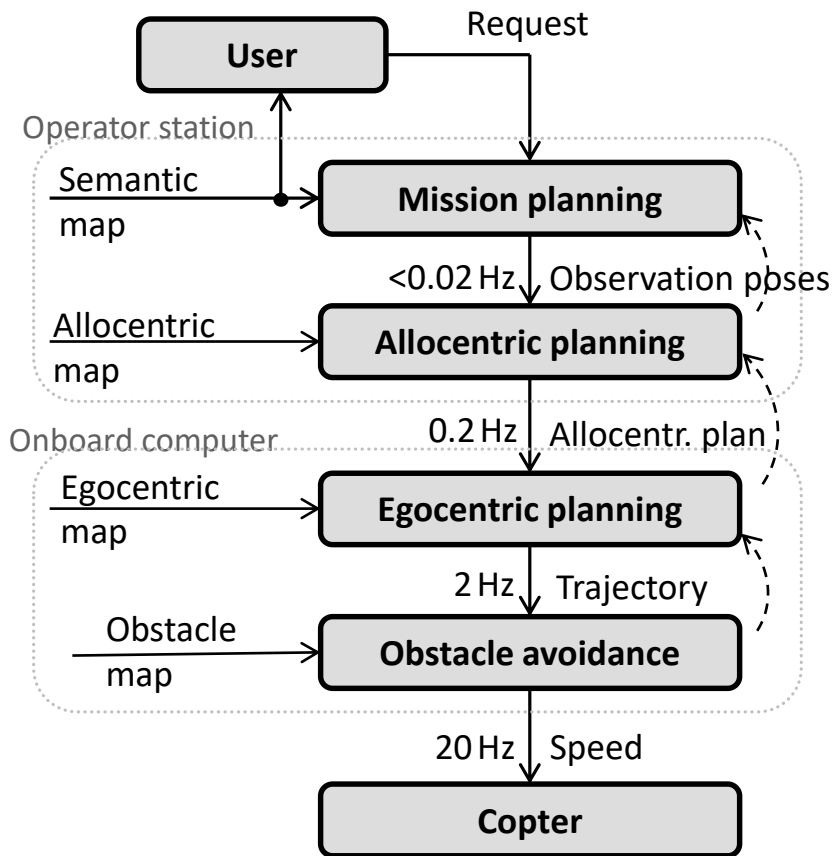
- Direct functional grasps not always feasible
- Pick up object with support hand, such that it can be grasped in a functional way



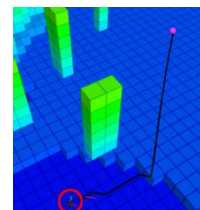
# Regrasping Experiments



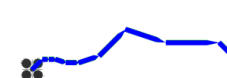
# Micro Aerial Vehicles: Hierarchical Navigation



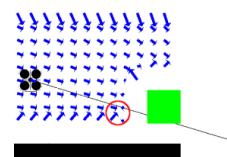
Mission plan



Allocentric planning

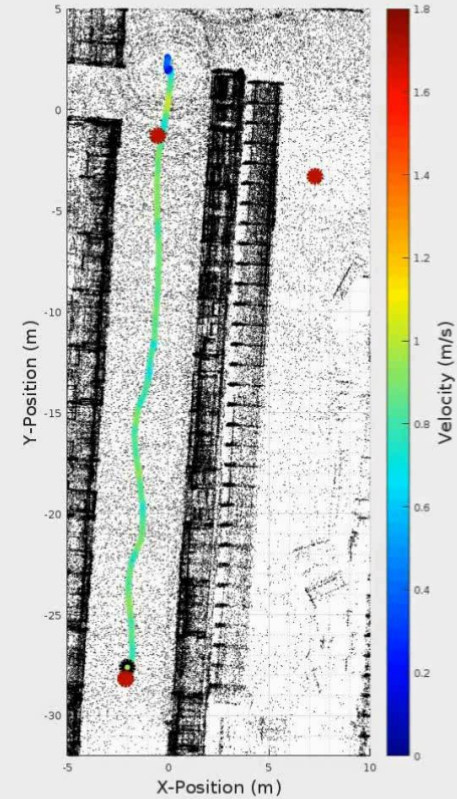
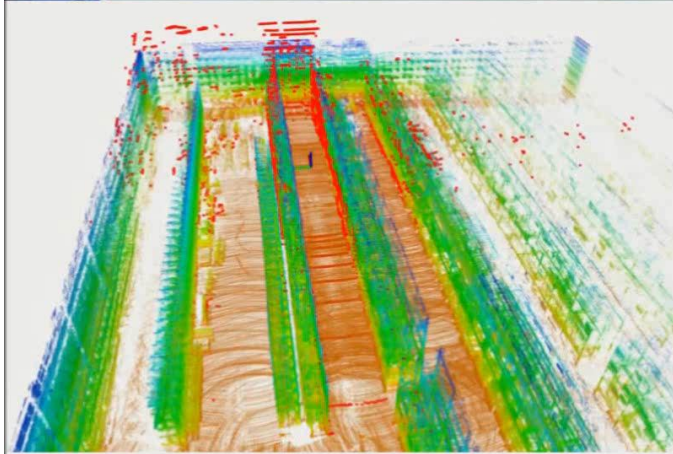


Egocentric planning



Obstacle avoidance

# InventAIRy: Autonomous Navigation in a Warehouse



# InventAIRy: Detected Tags in Shelf



## Initial demonstrator



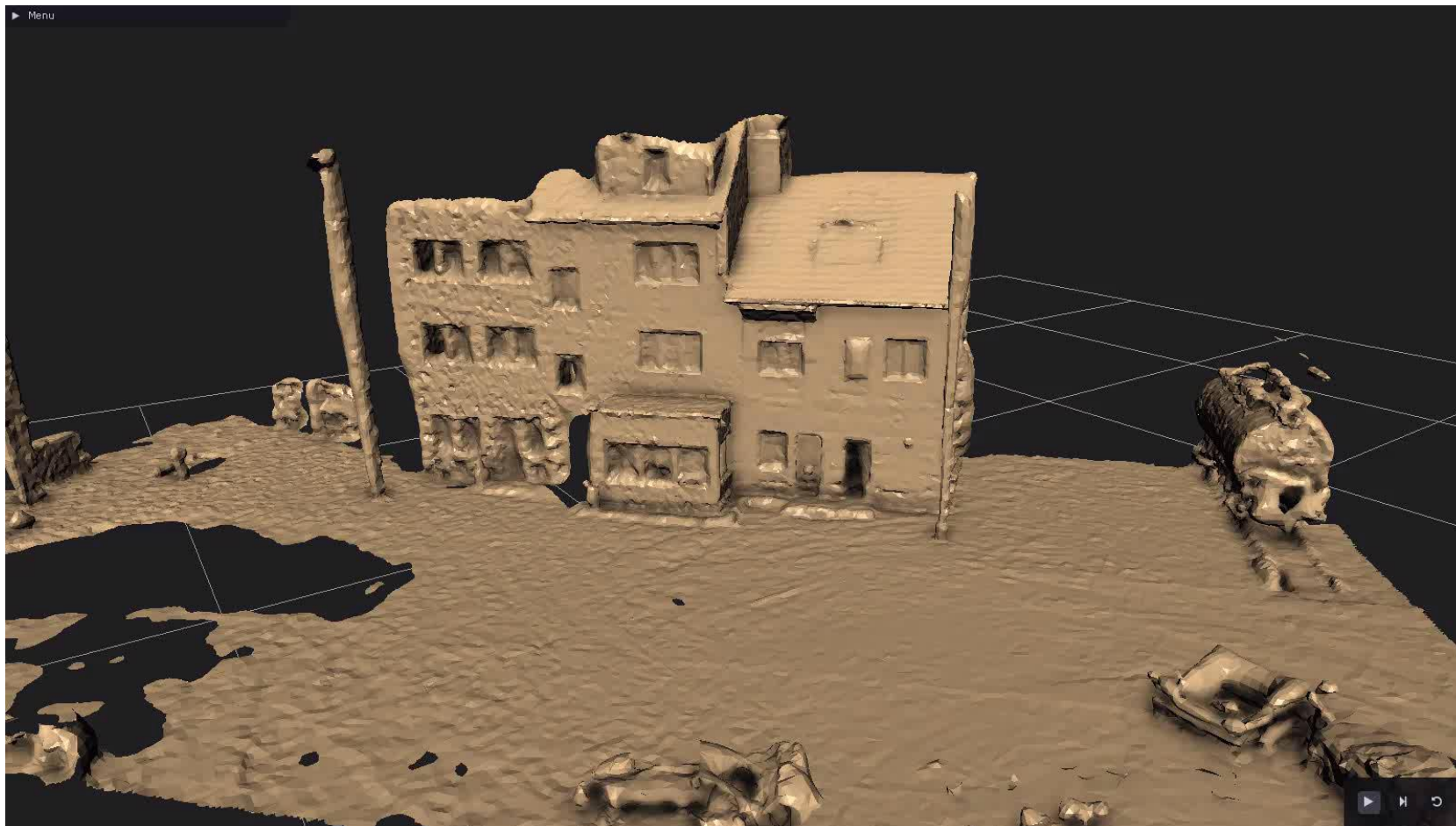
- Basis: DJI Matrice 600 Pro
- Sensors: Velodyne VLP 16, FLIR Boson, 2x FLIR BlackFly S
- Tilttable sensor head

## Current demonstrator



- Basis: DJI Matrice 210 v2
- Sensors: Ouster OS-0, FLIR AGX, 2x Intel RealSense D455
- IP43 water resistance

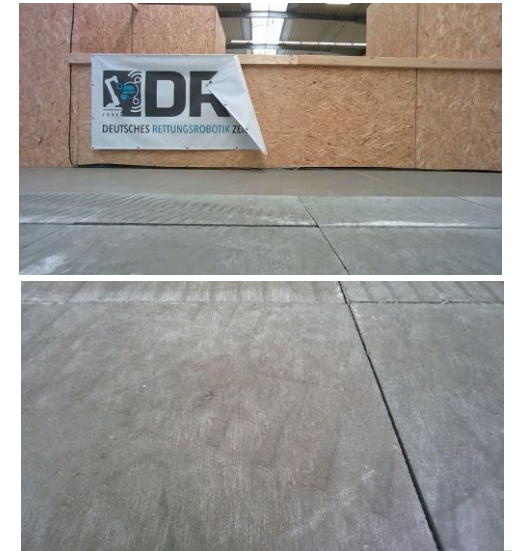
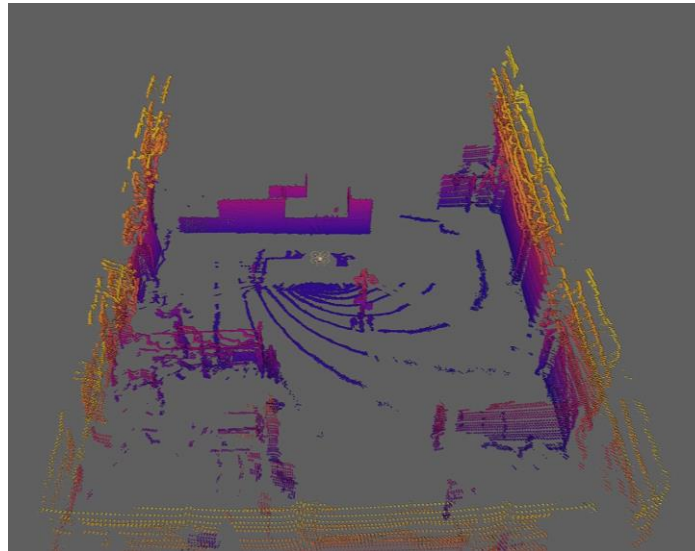
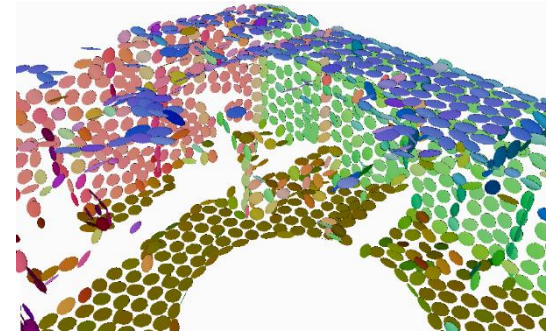
# Modeling the Brandhaus Dortmund





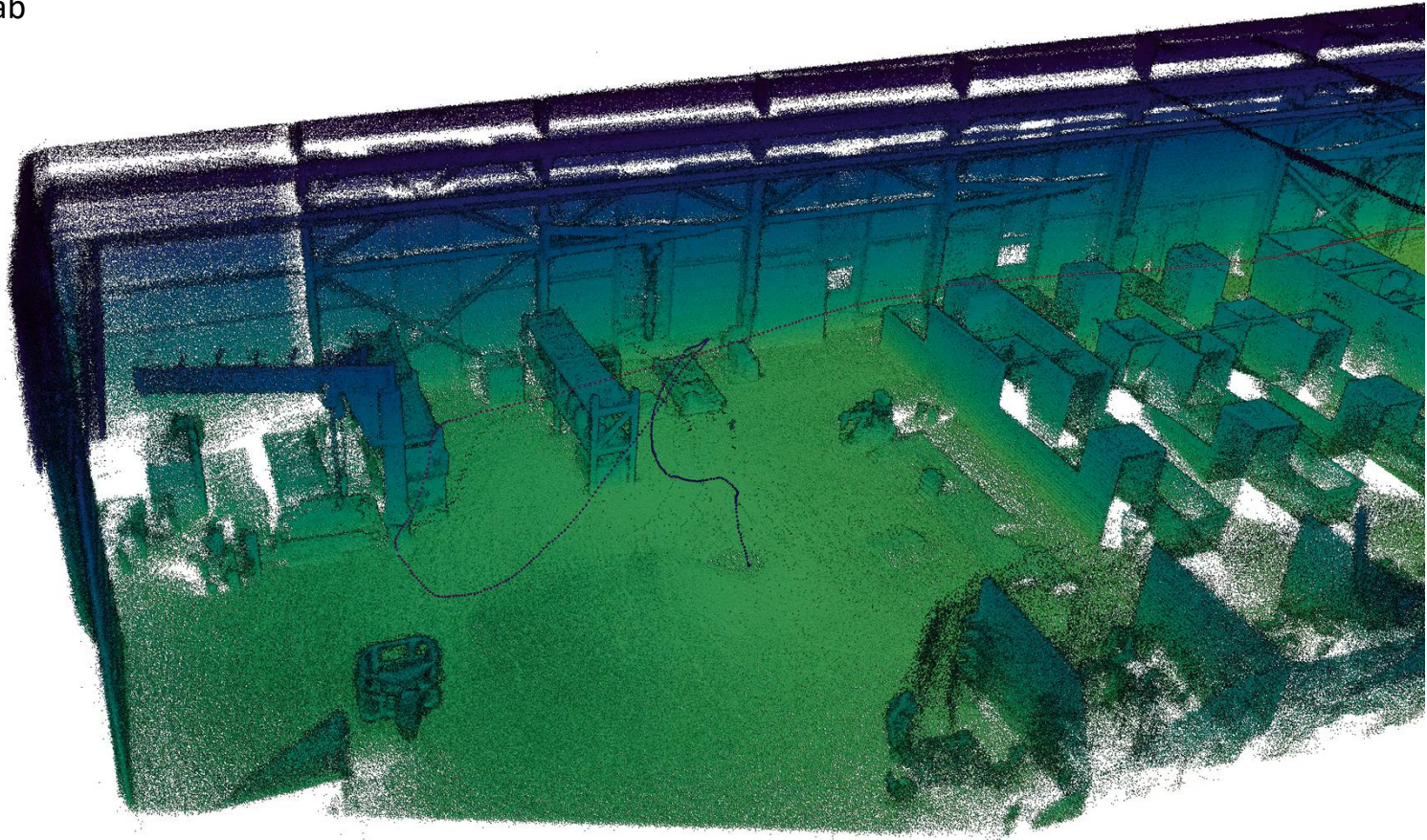
# Real-time LiDAR Odometry with Continuous-time Trajectory Optimization

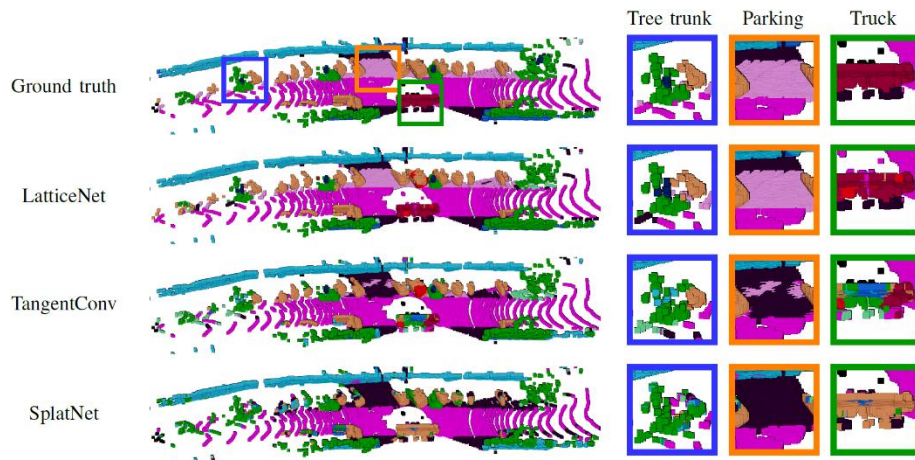
- Simultaneous registration of multiple multiresolution surfel maps using Gaussian mixture models and temporally continuous B-spline
- Accelerated by sparse permutohedral voxel grids and adaptive choice of resolution
- Real-time onboard processing 16-20 Hz
- Open-Source  
[https://github.com/AIS-Bonn/lidar\\_mars\\_registration](https://github.com/AIS-Bonn/lidar_mars_registration)



# 3D LiDAR Mapping

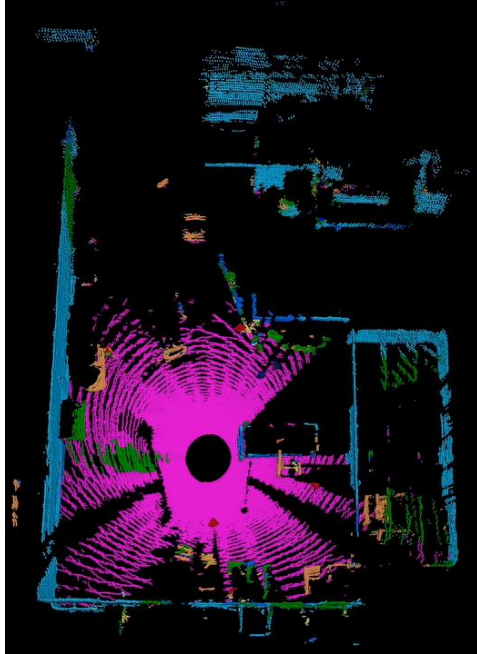
DRZ Living Lab





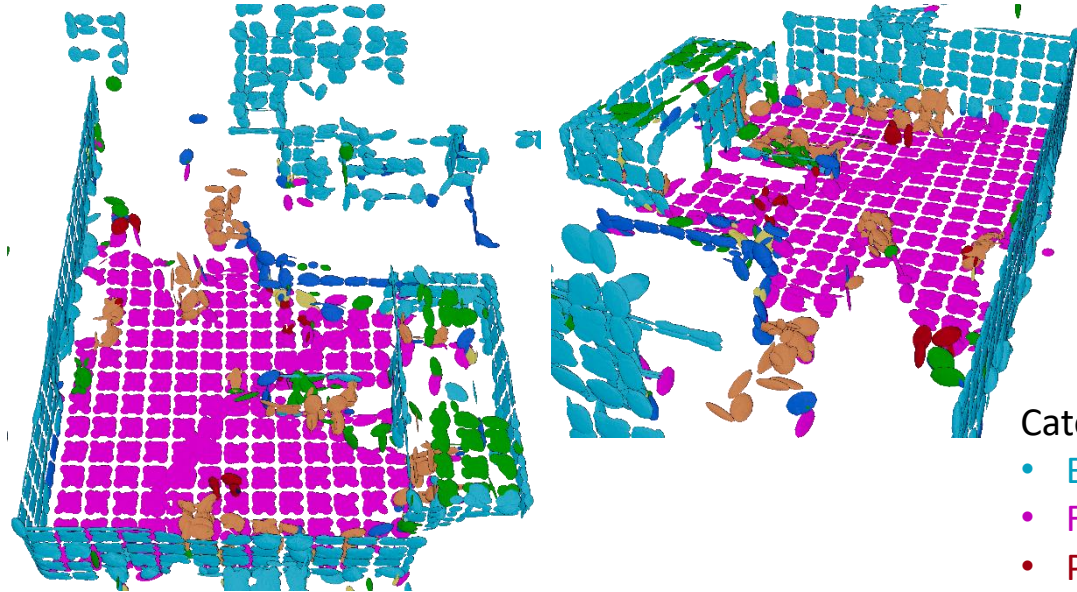
- LatticeNet segmentation of 3D point clouds based on sparse permutohedral grid
- Hierarchical information aggregation through U-Net architecture
- LatticeNet is real-time capable and achieves excellent results in benchmarks

# Semantic Fusion: 3D LiDAR Mapping



Segmented point cloud

Minimax-Viking fire house



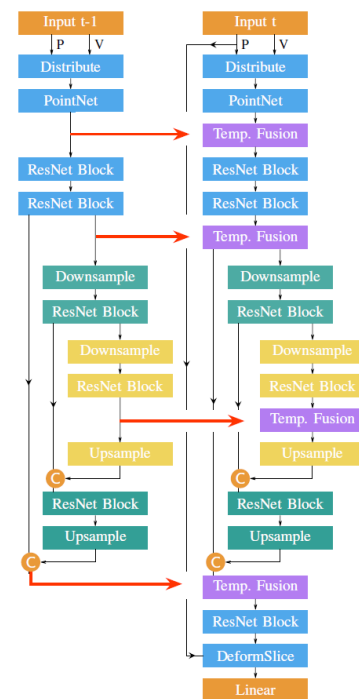
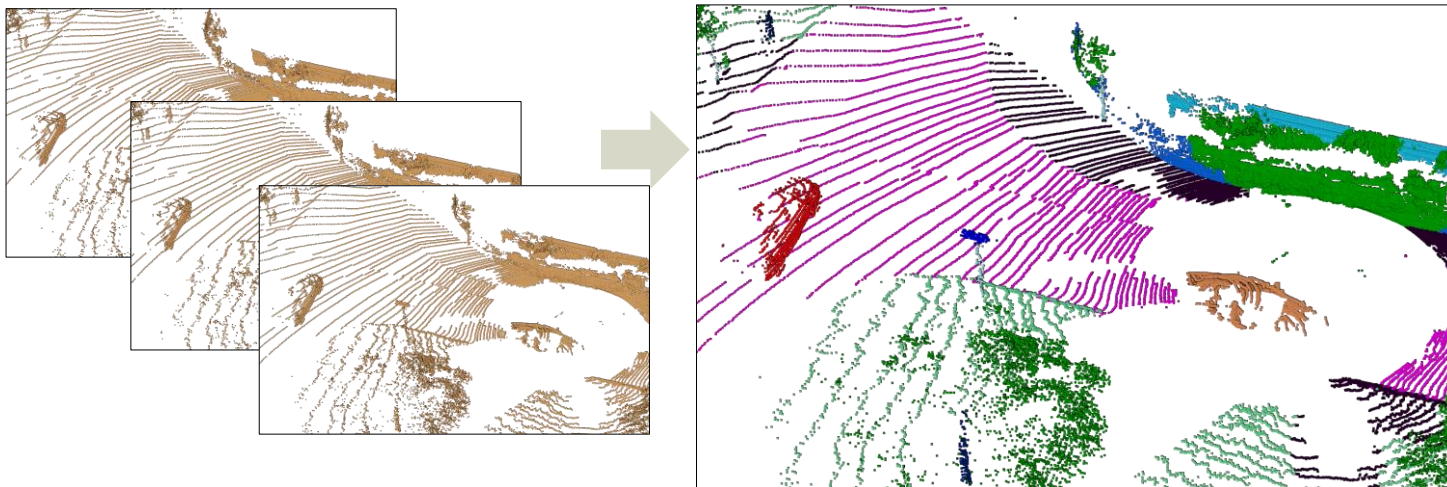
Semantic multiresolution surfel map

Categories:

- Building
- Floor
- Persons
- Vehicles
- Fence
- Vegetation

# Semantic Fusion: Temporal LatticeNet

- Semantic segmentation of sequences of 3D point clouds
- Integration of recurrent connections
- Trained on three scans of SemanticKITTI
- Distinguishing moving from parking vehicles

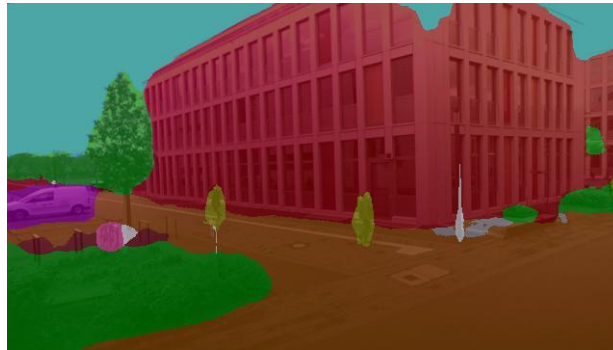
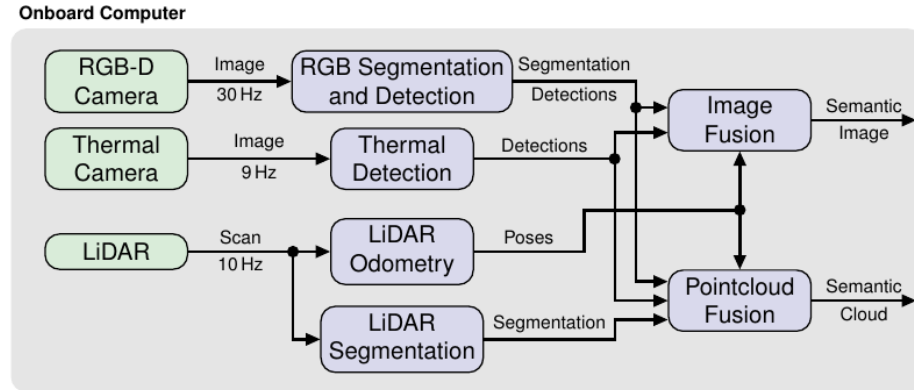


Categories:

- Street
- Moving Vehicle
- Parking Vehicle
- Vegetation

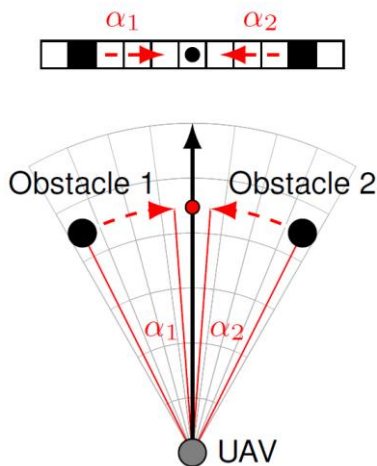
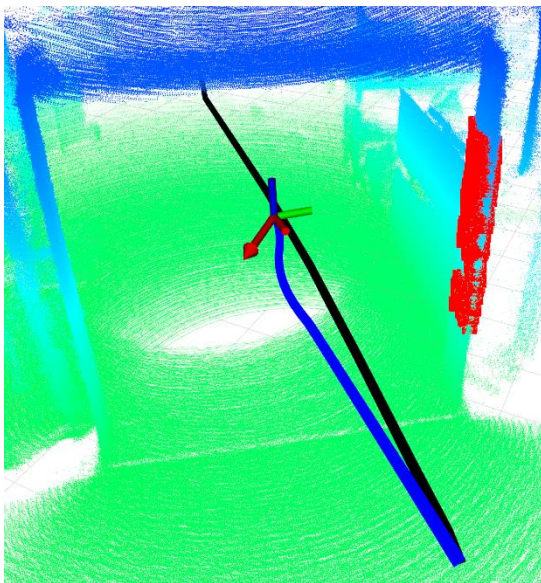
# Onboard Multimodal Semantic Fusion

- Real-time semantic segmentation and object detection ( $\approx 9\text{Hz}$ ) with EdgeTPU / iGPU
  - SalsaNext for LiDAR
  - DeepLabv3 for RGB images
  - SSD MobileDet for Thermal/RGB
- Late-fusion for
  - Point cloud
  - Image segmentation

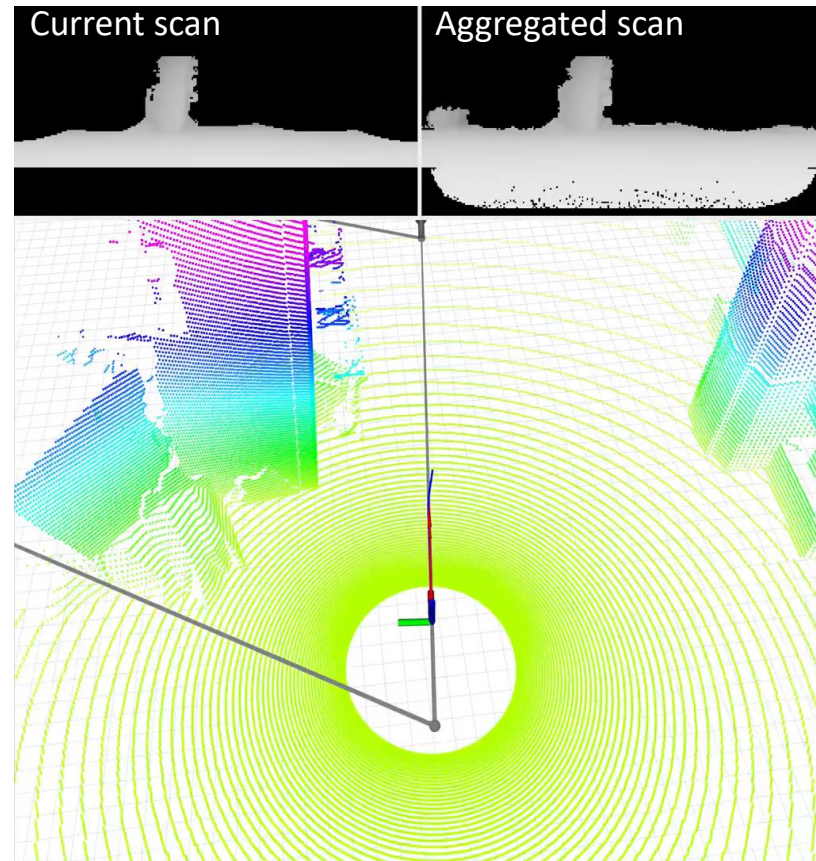


# Predictive Angular Potential Field-based Obstacle Avoidance

- Aggregate LiDAR scans in range image
- Adjust direction using angular potential field
- Predict trajectory and range image
- Scale velocity based on time-to-contact

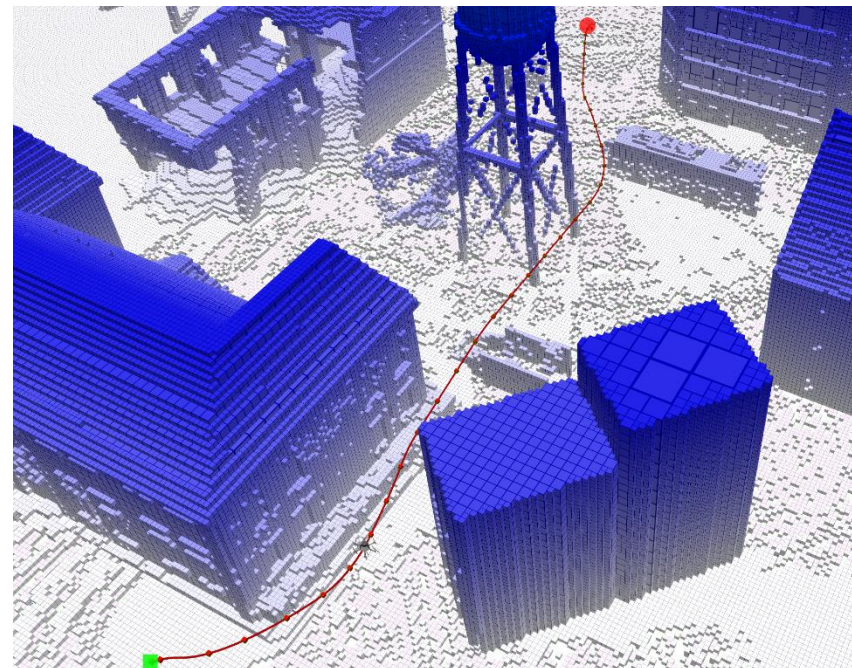
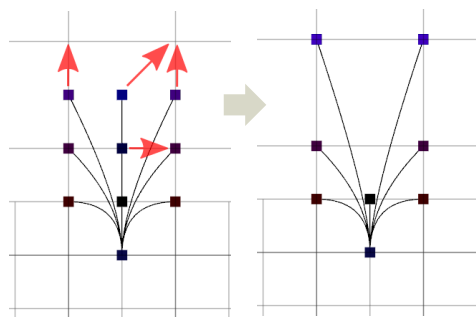
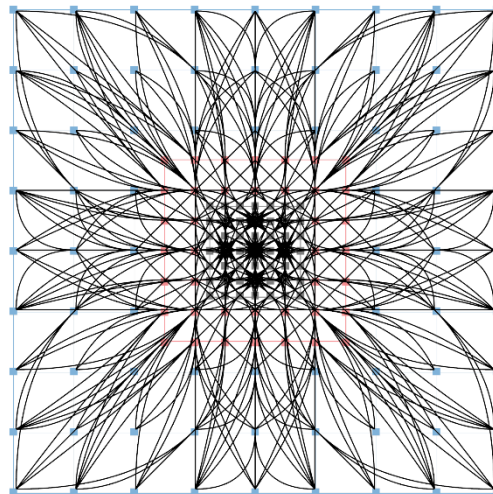


Angular Potential Field



# Dynamic 3D Navigation Planning

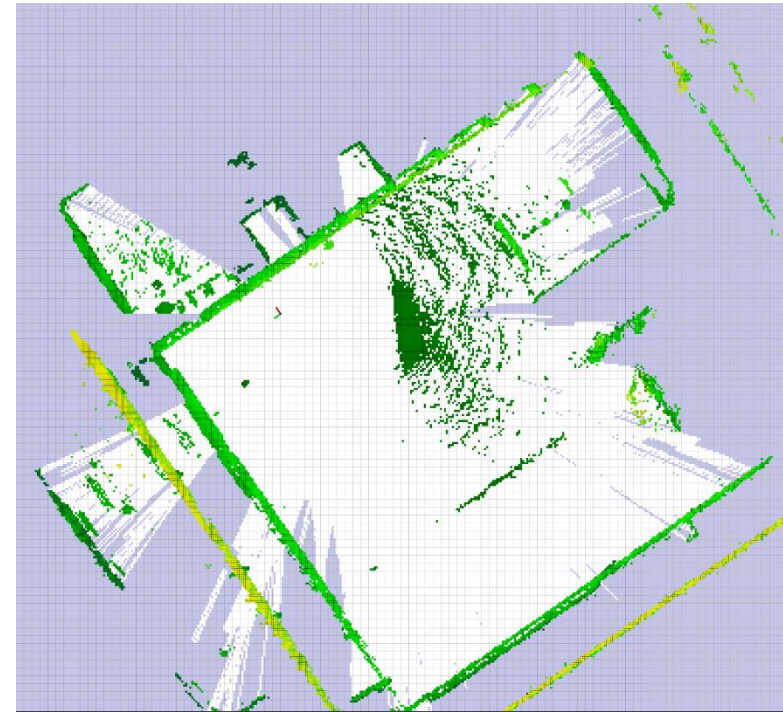
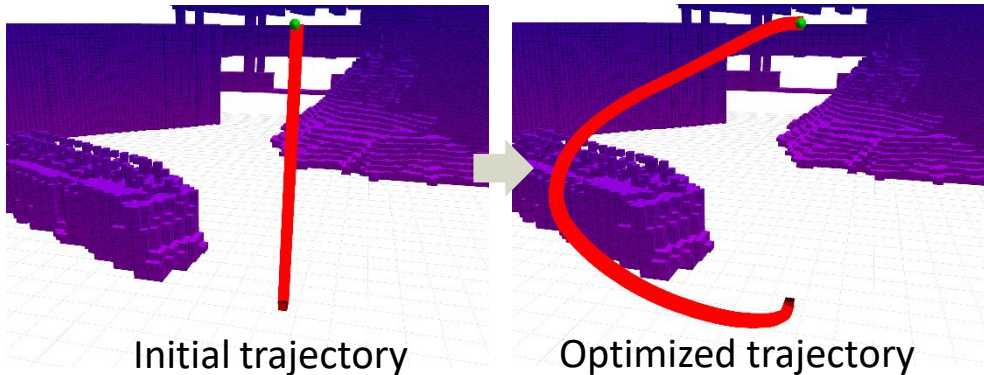
- Positions and velocities in sparse local multiresolution grid
- Adaptation of movement primitives to grid
- Optimization of flight time and control costs
- 1 Hz replanning





# Planning with Visibility Constraints

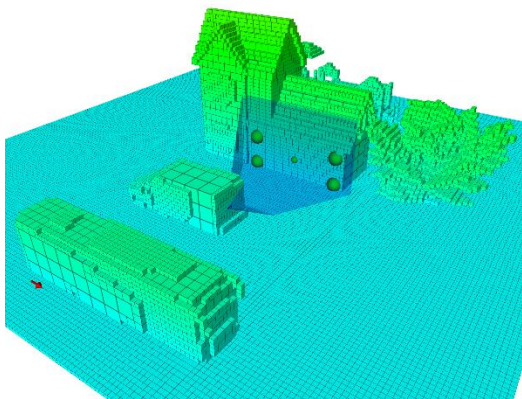
- Extra costs for flight through unmapped volumes
- Consideration of sensor frustum:
  - Coupling of vertical and horizontal motion
  - Preferred forward flight with limited rotational speed



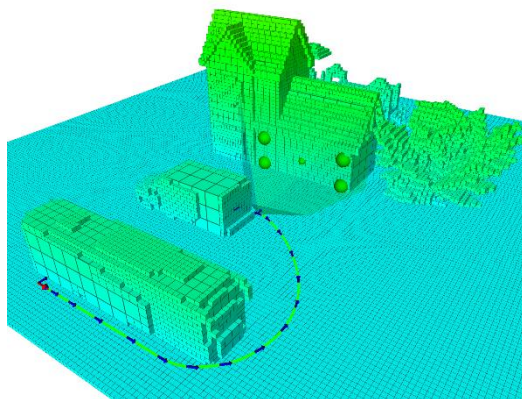
Obstacle map

# Observation Pose Planning

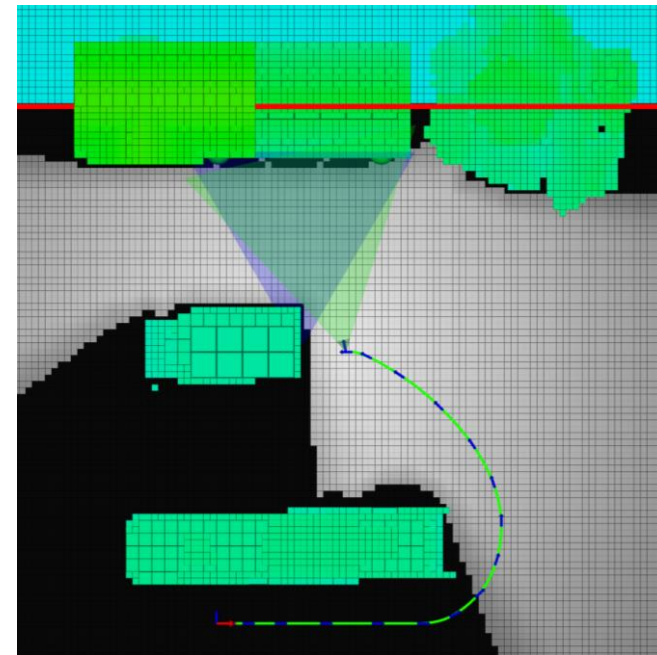
- Planning of observation poses with line of sight to the target object despite occlusions
- Target objects are defined by position, line of sight and distance
- Optimization of observation poses with regard to visibility quality and accessibility



Initial observation pose

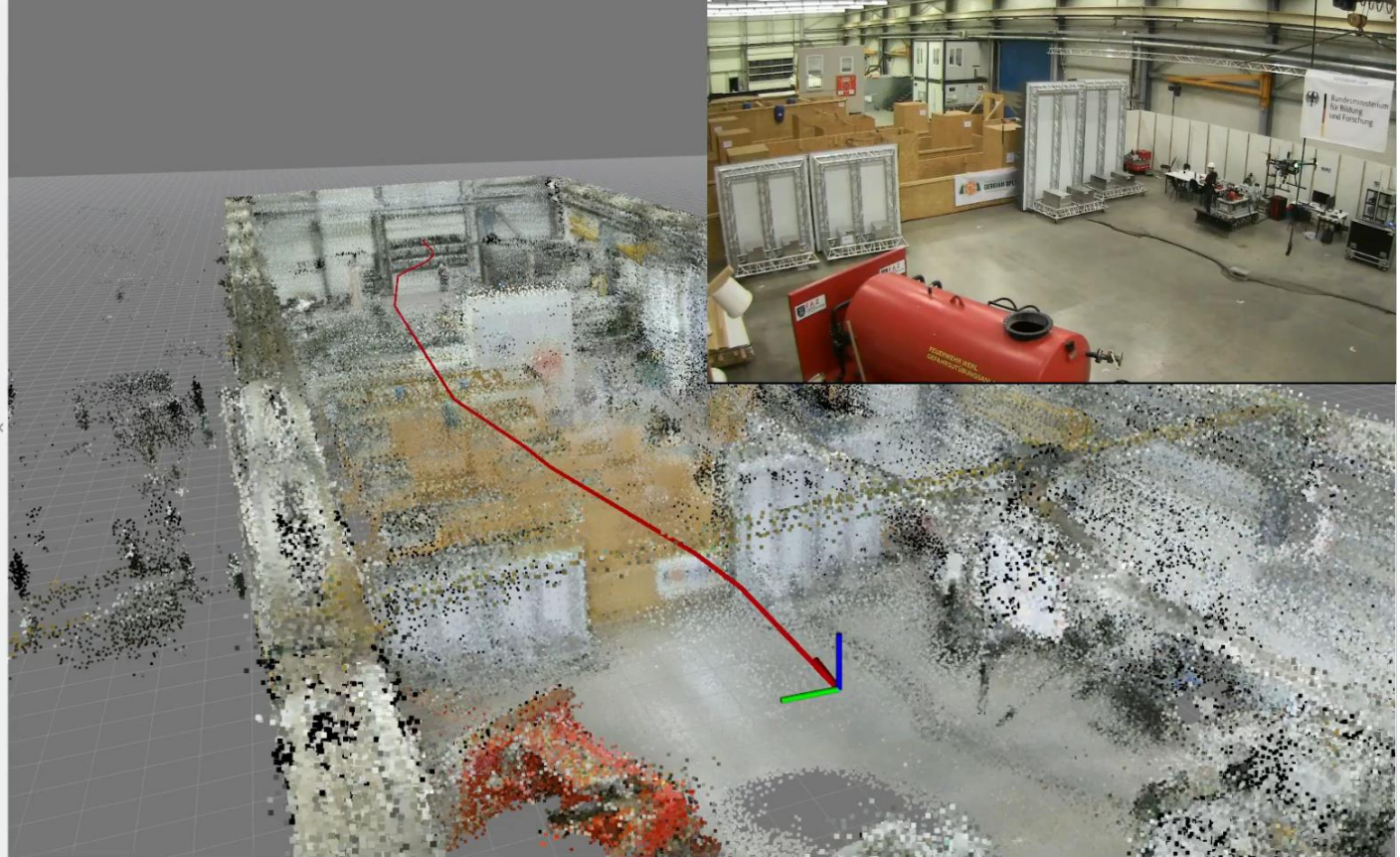
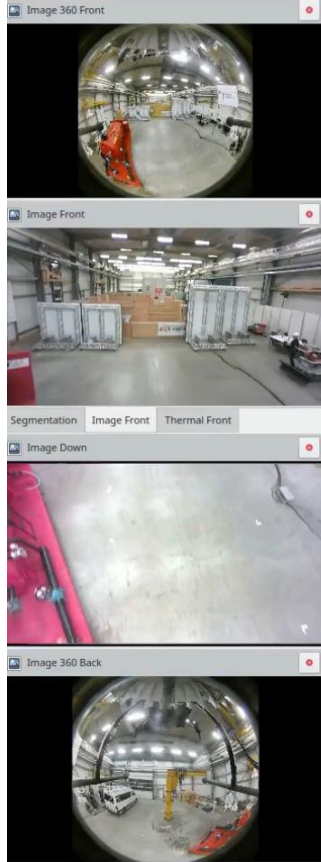


Optimized path



Top-down view

# Autonomous Flight without GNSS

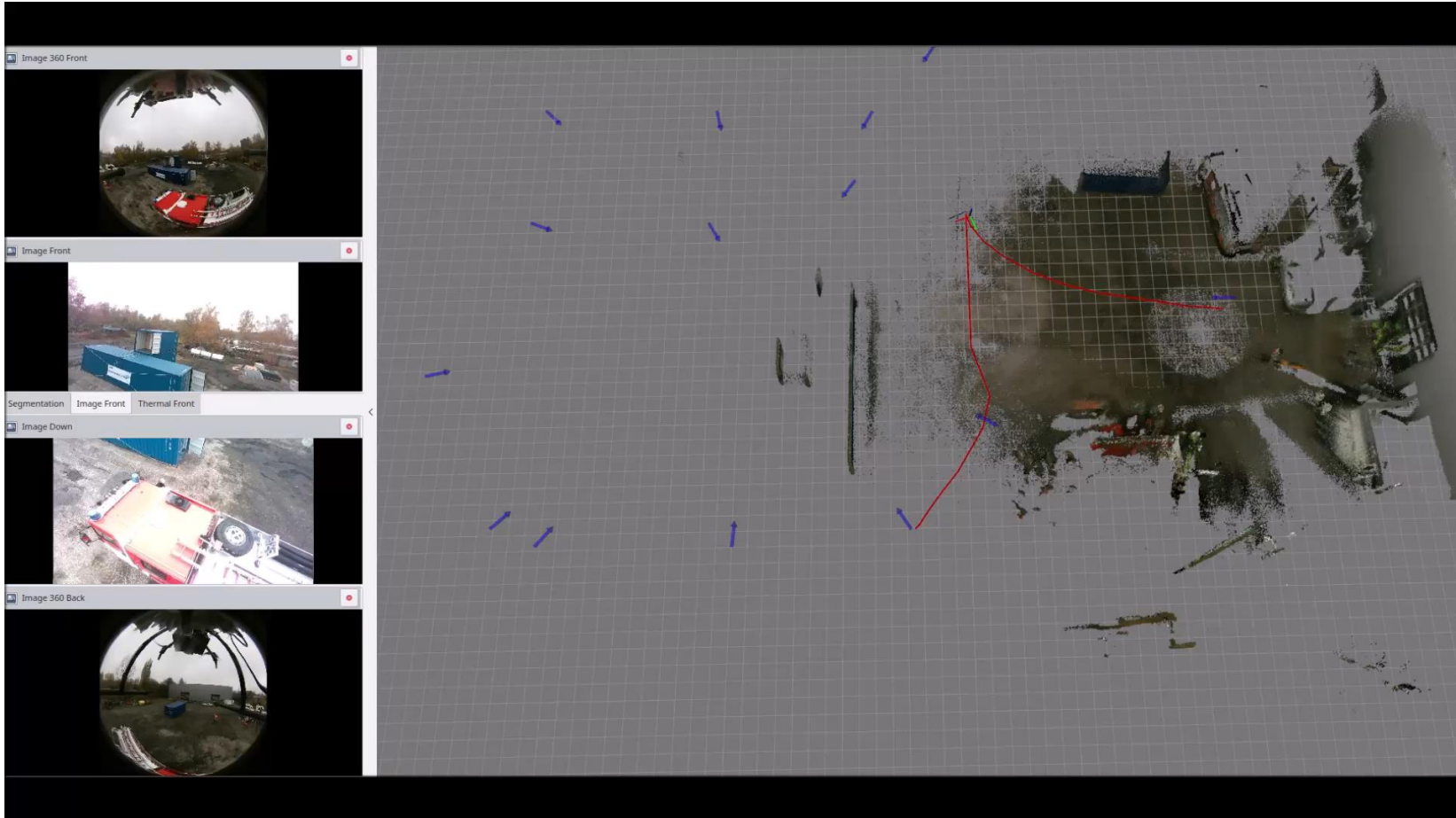


# Exploration

- Definition of target area w.r.t. satellite images or maps
- Simple exploration patterns (spirals, meanders, ...)
- Collision check
- TSP to determine segment sequence
- Continuous replanning



Campus Poppelsdorf

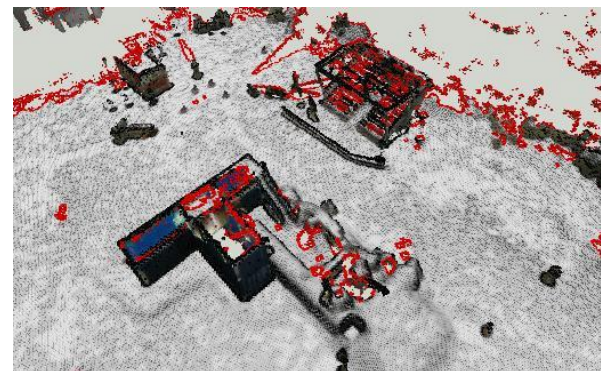


# Terrain Classification for Traversability

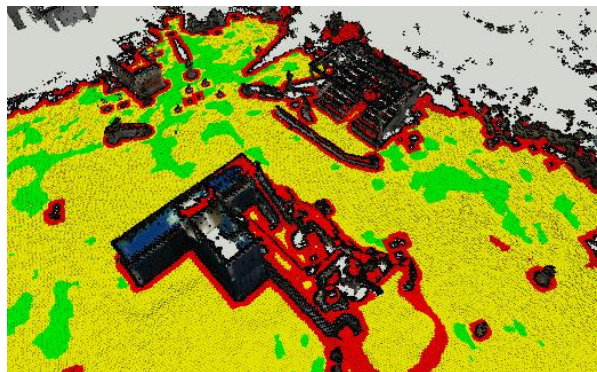
- Based on voxel-filtered aggregated point cloud
- Terrain classification based on local height differences in the robot ground robot footprints
- Categories: drivable, walkable, unpassable
- Reachability analysis



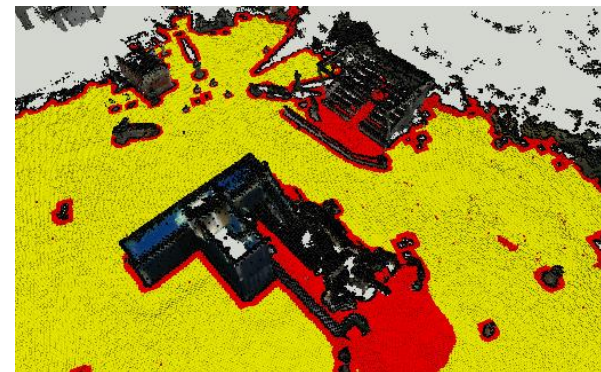
Aggregated colored point cloud



Local height differences



Terrain category

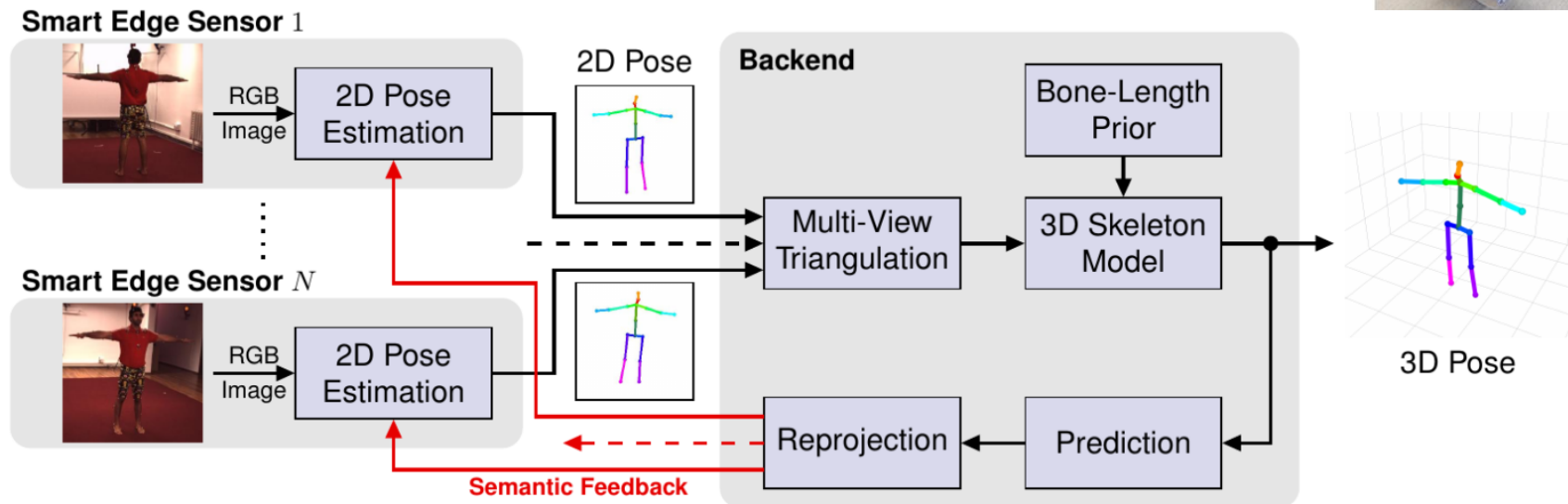


Reachability

[Schleich et al., ICUAS 2021]

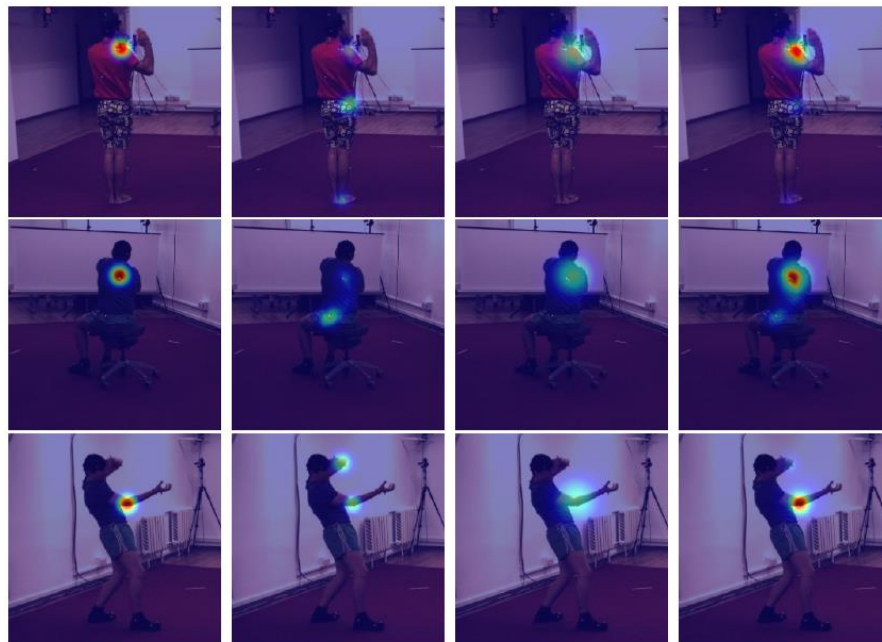
# Real-Time Multi-View 3D Human Pose Estimation using Semantic Feedback to Smart Edge Sensors

- Triangulation and skeleton model to recover 3D pose
- Semantic feedback channel for bidirectional communication between backend and sensors



# Real-Time Multi-View 3D Human Pose Estimation using Semantic Feedback to Smart Edge Sensors

- Feedback heatmap is rendered from feedback skeleton and fused with detection on sensors
- Feedback heatmap helps to recover from incorrect or imprecise 2D joint detections
- Examples:
  - Occluded left wrist (rows 1 and 2)
  - Confusion of left and right elbow (row 3)

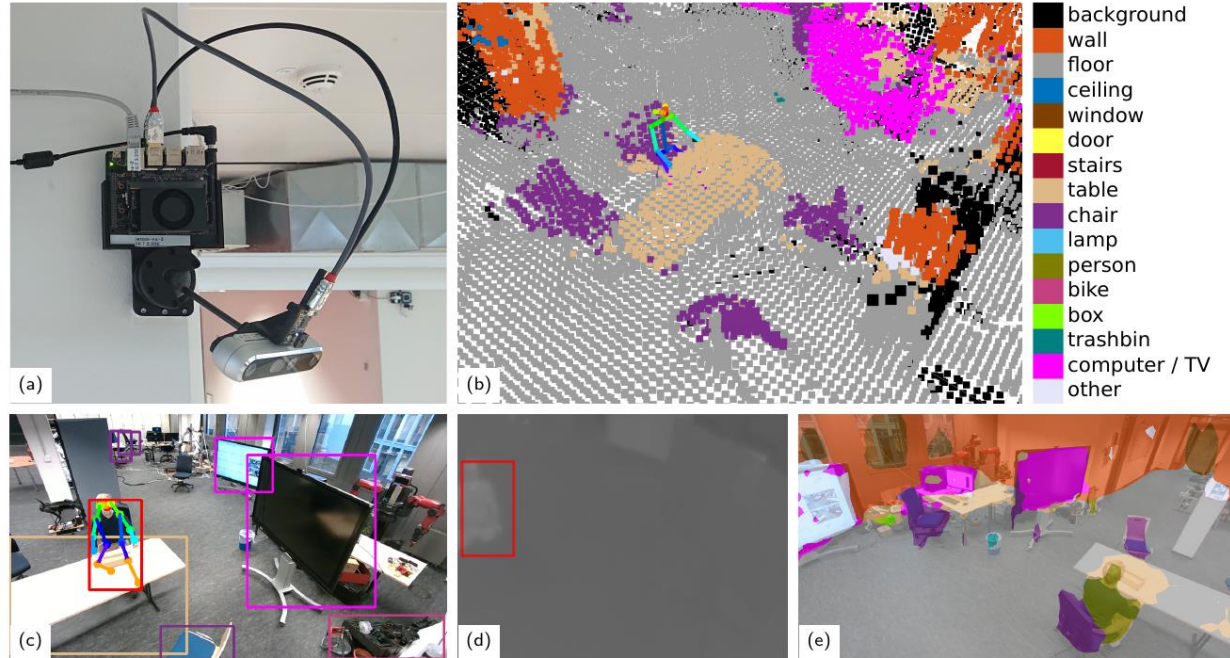
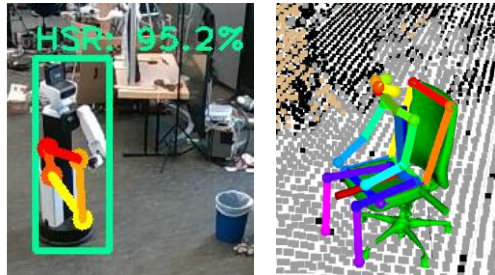


(a) ground-truth (b) detected (c) feedback (d) fused



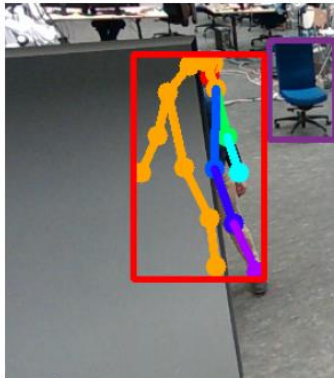
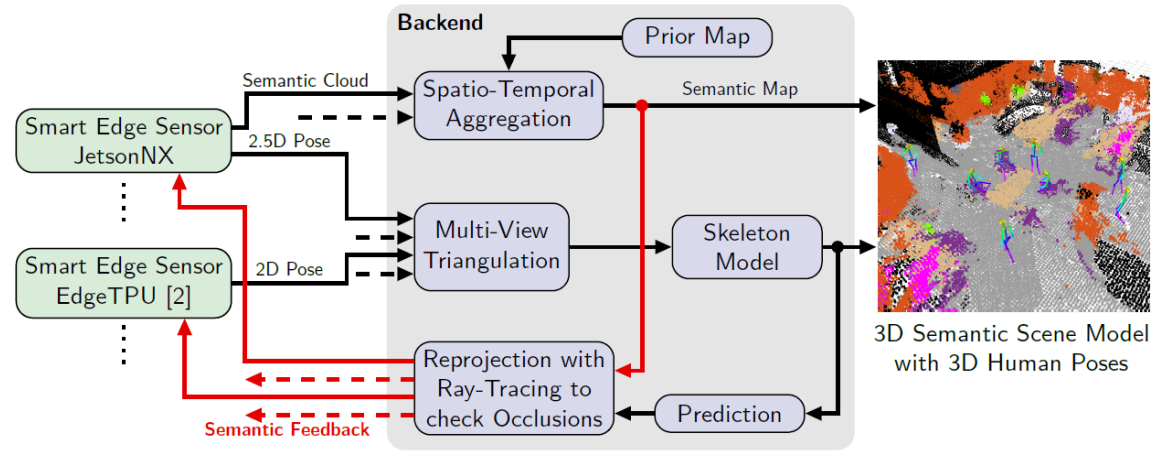
# Semantic Perception with Smart Edge Sensor Network

- Object detection and semantic segmentation of RGB images
- Person detection in IR images
- Semantic labelling of RGB-D point clouds
- Pose estimation for mobile robot and objects

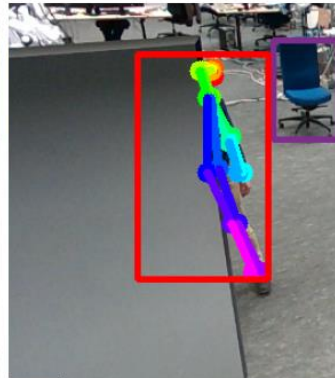


# 3D Human Pose Estimation with Occlusion Feedback

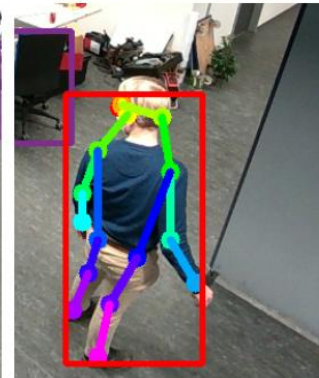
- Heavy occlusion causes the pose estimation to collapse to the visible side only
- With occlusion feedback occluded joint detections can be discarded and the local model is completed



With occlusion feedback



W/o occlusion feedback



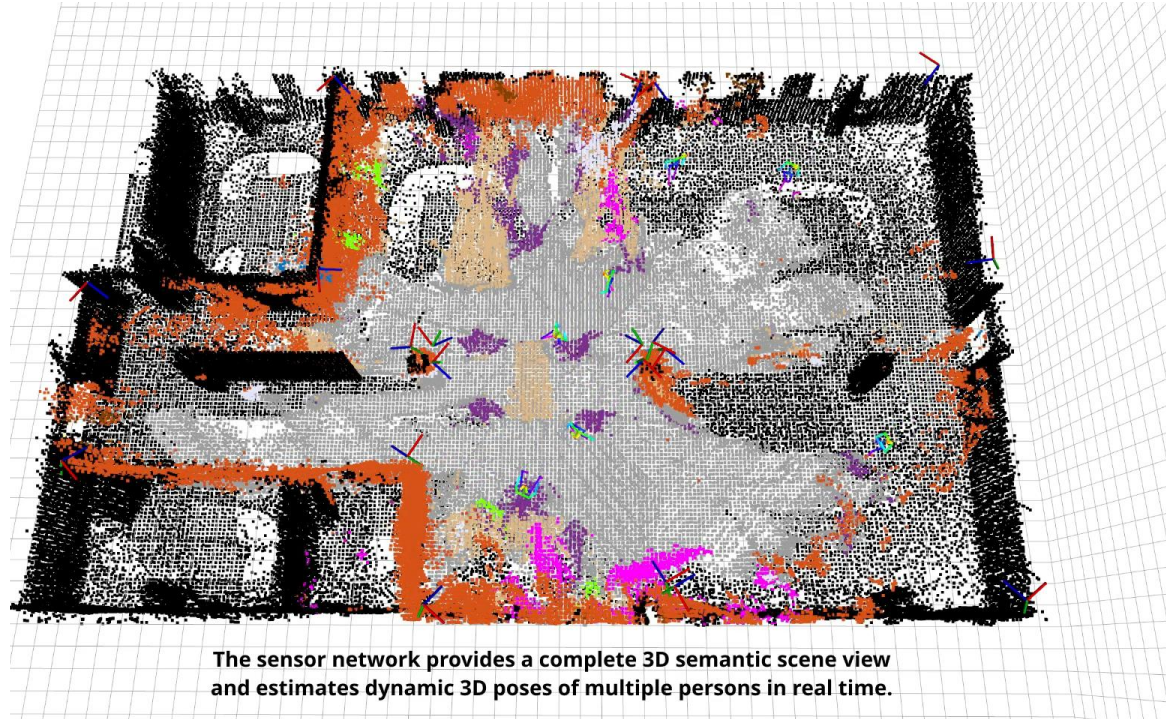
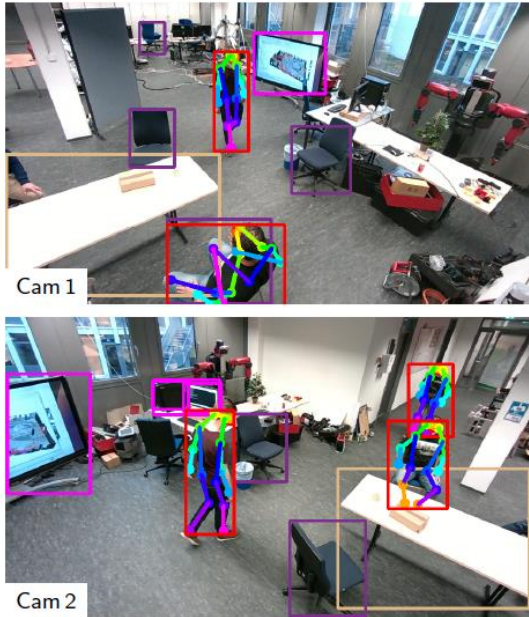
Unoccluded reference



Fully occluded

# Evaluation in Real-World Multi-Person Scenes

- 20 smart edge sensors (4 Jetson NX, 16 Edge TPU), covering 12×22 m area
- Experiments with 8 persons moving through the scene

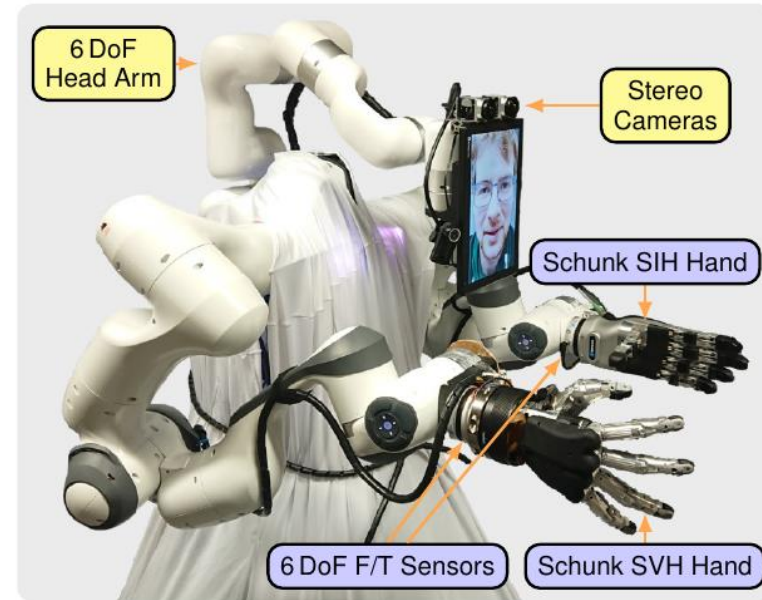
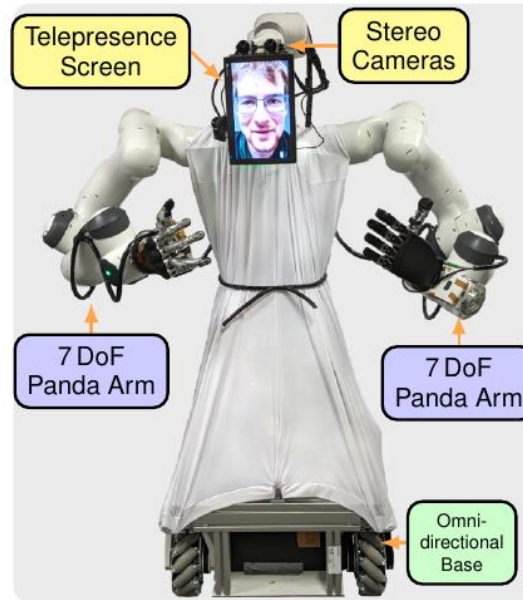
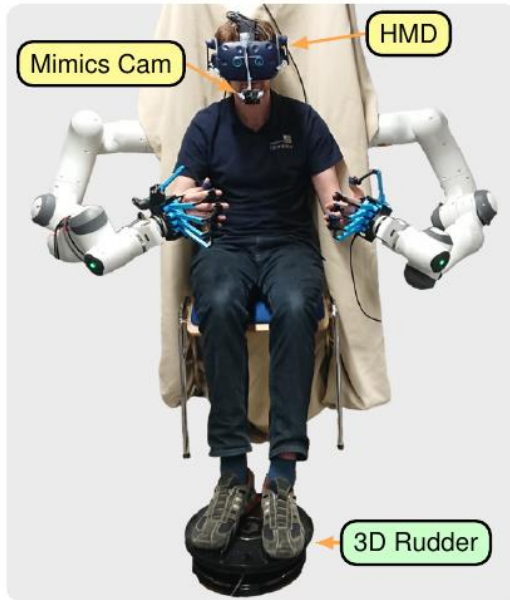


# \$10M ANA Avatar XPRIZE Competition

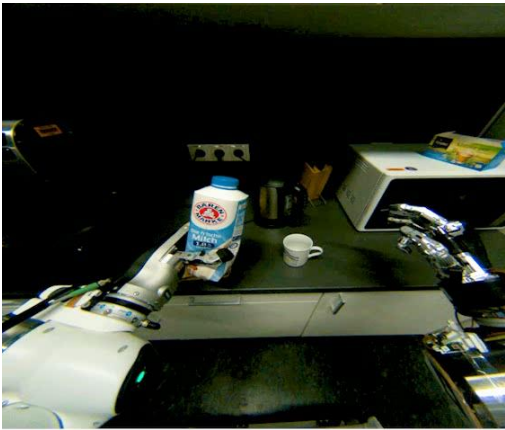
- Requires mobility, manipulation, human-human interaction
- Focuses on the immersion in the remote environment and the presence of the remote operator



- Two-armed avatar robot designed for teleoperation with immersive visualization & force feedback
- Operator station with HMD, exoskeleton and locomotion interface



# Team NimbRo Semifinal Submission



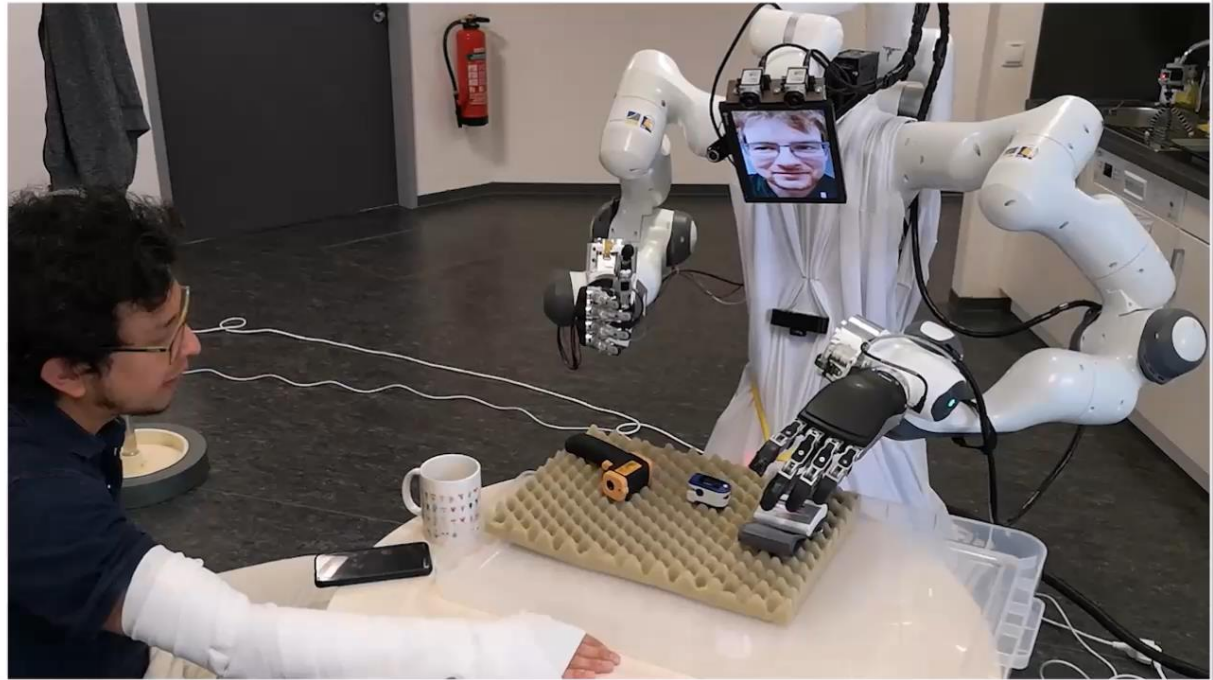
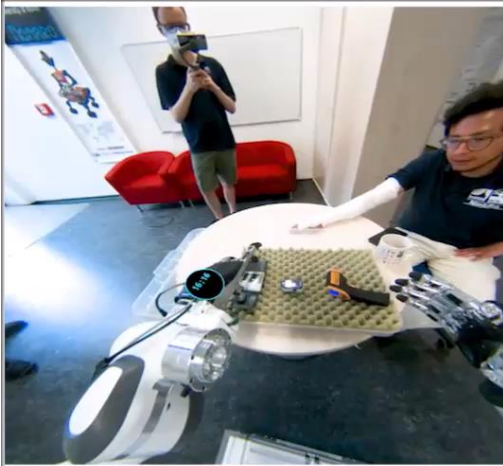
# Team NimbRo

## Semifinal Team Video

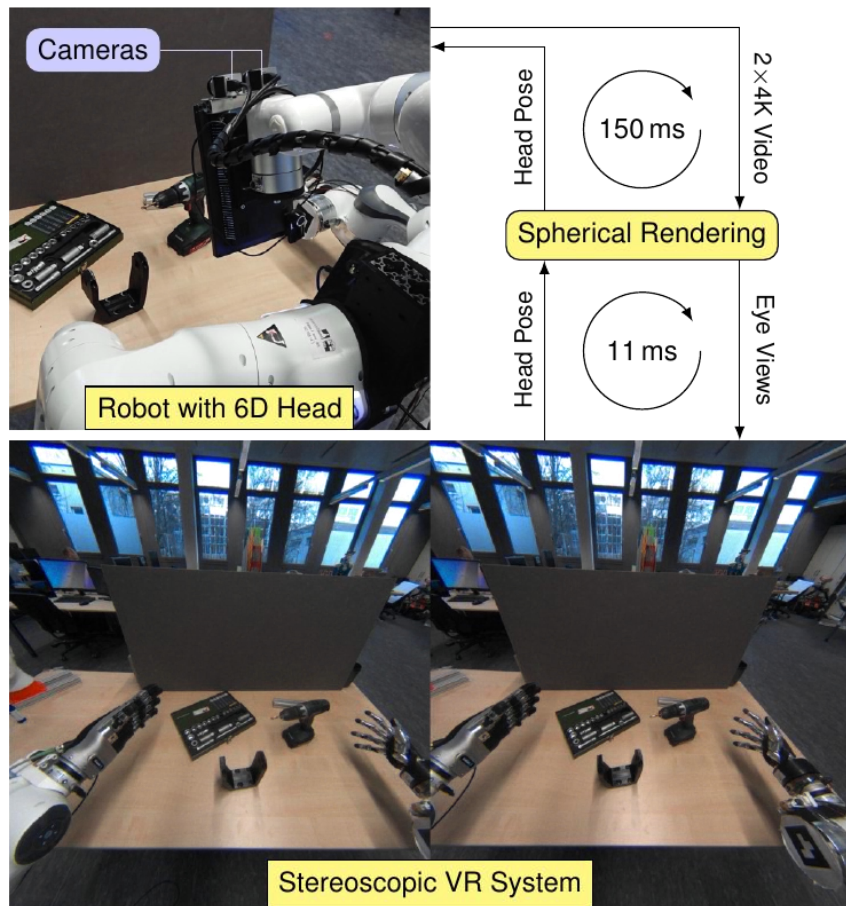


### Tasks

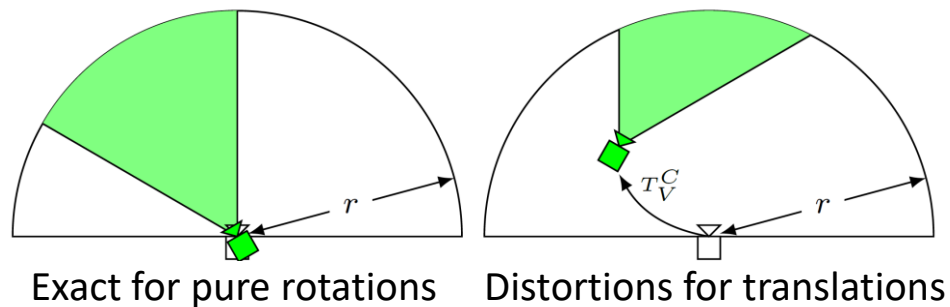
1. Make a coffee
2. Greet the recipient
3. Measure temperature
4. Measure blood pressure
5. Measure oxygen saturation
6. Help recipient with jacket



# NimbRo Avatar: Immersive Visualization



- 4K wide-angle stereo video stream
- 6D neck allows full head movement
  - Very immersive
- Spherical rendering technique hides movement latencies
  - Assumes constant depth





# NimRo Avatar: Operator Face Animation

- Operator images without HMD
- Capture mouth and eyes
- Estimate gaze direction and facial keypoints
- Generate animated operator face using a warping neural network



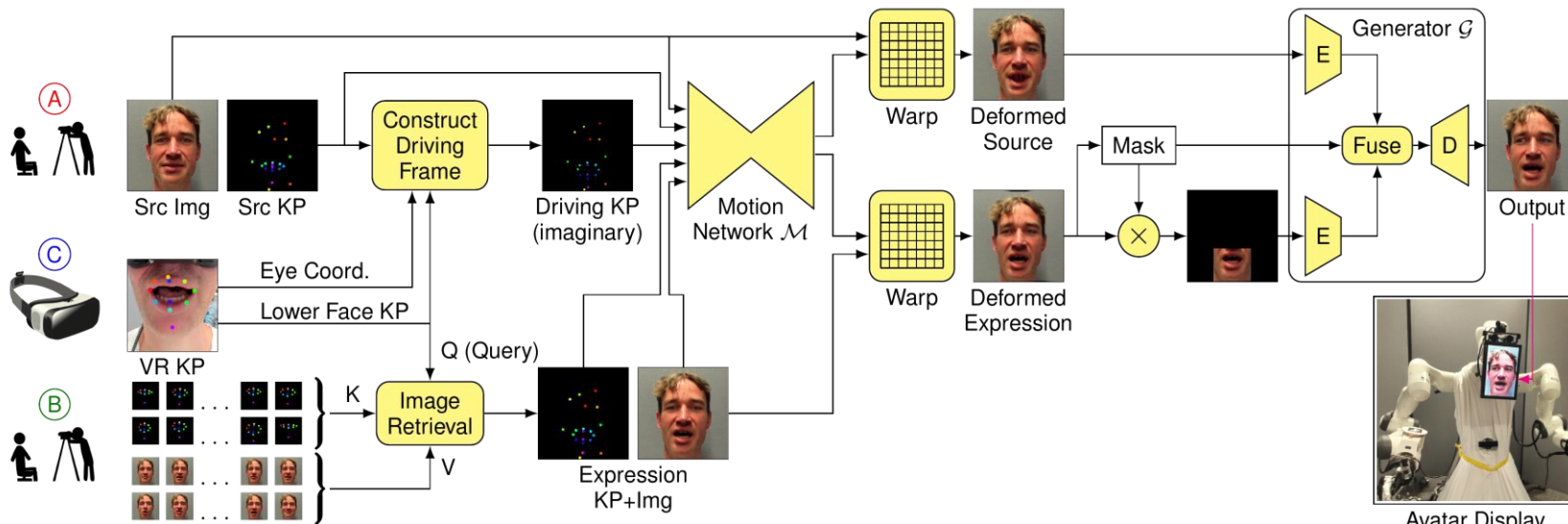
Left Eye



Mouth



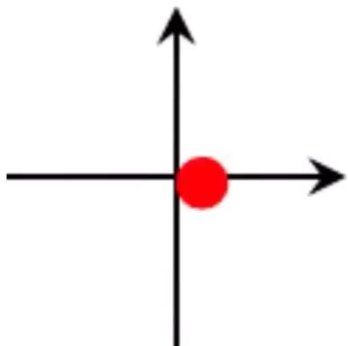
Right Eye



[Rochow et al. IROS 2022]

# NimbRo Avatar: Operator Face Animation

Gaze  
Direction



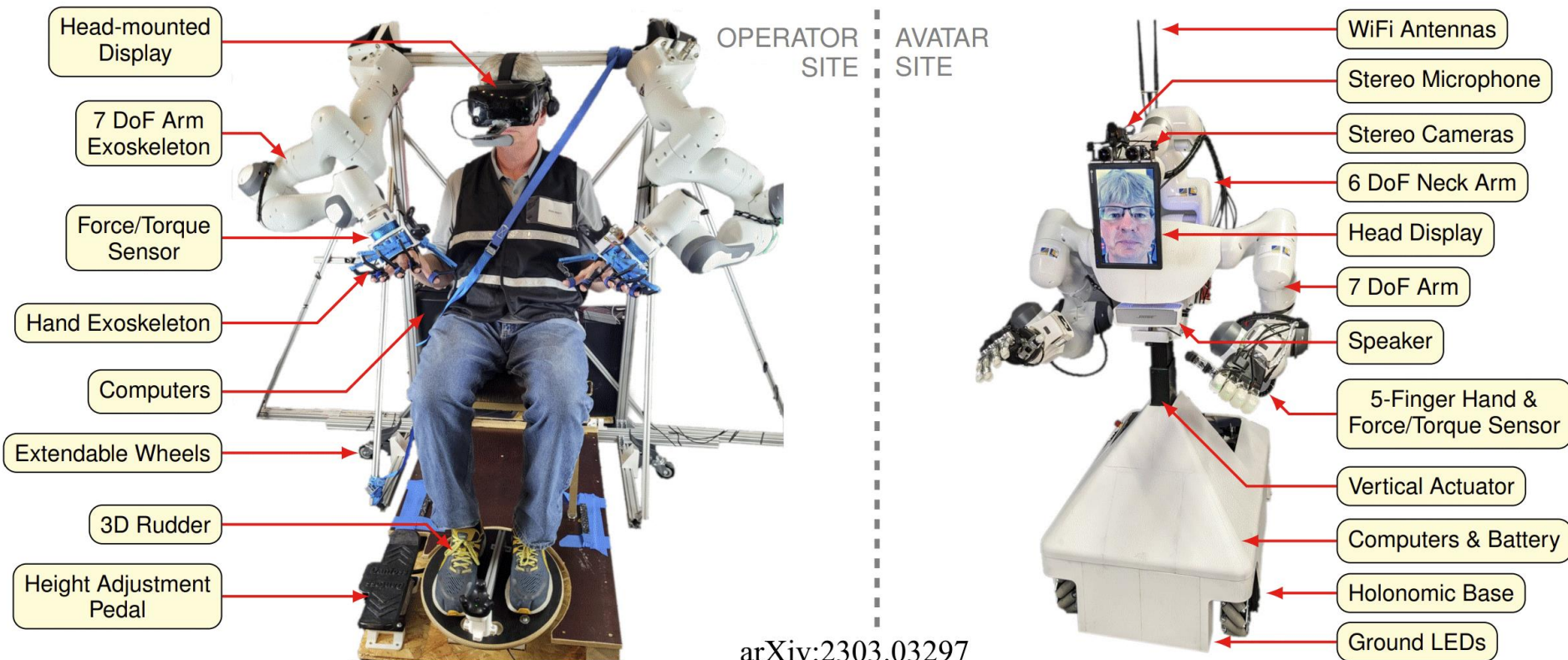
Output

Mouth Cam



# NimbRo Avatar System for ANA Avatar XPRIZE Finals

- New requirements: tetherless, remote perception of haptics, reliability



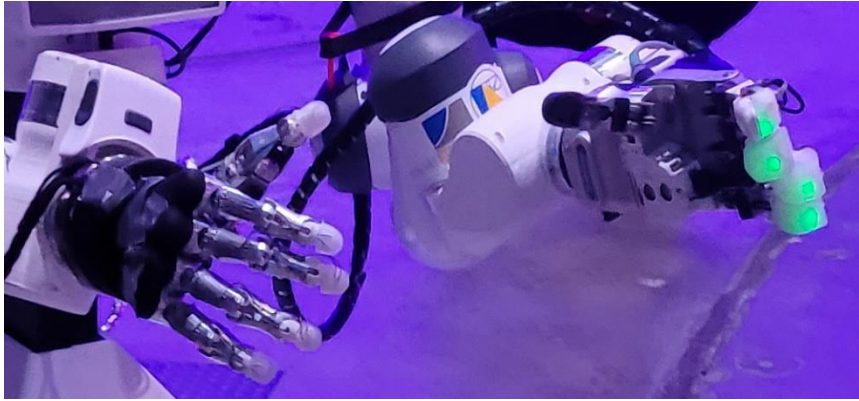
arXiv:2303.03297

[Schwarz, Lenz et al.]

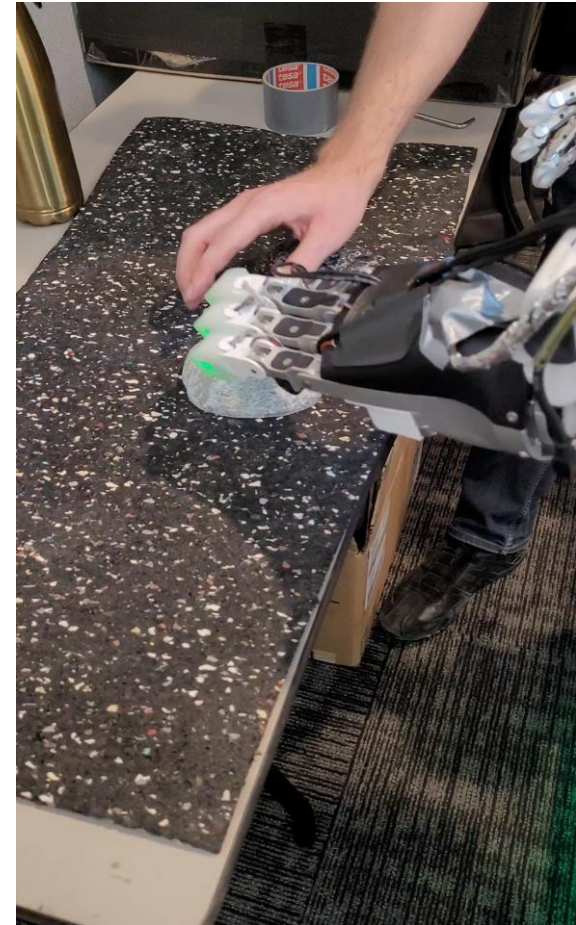


# Haptic Perception

- Sensors in the finger tips

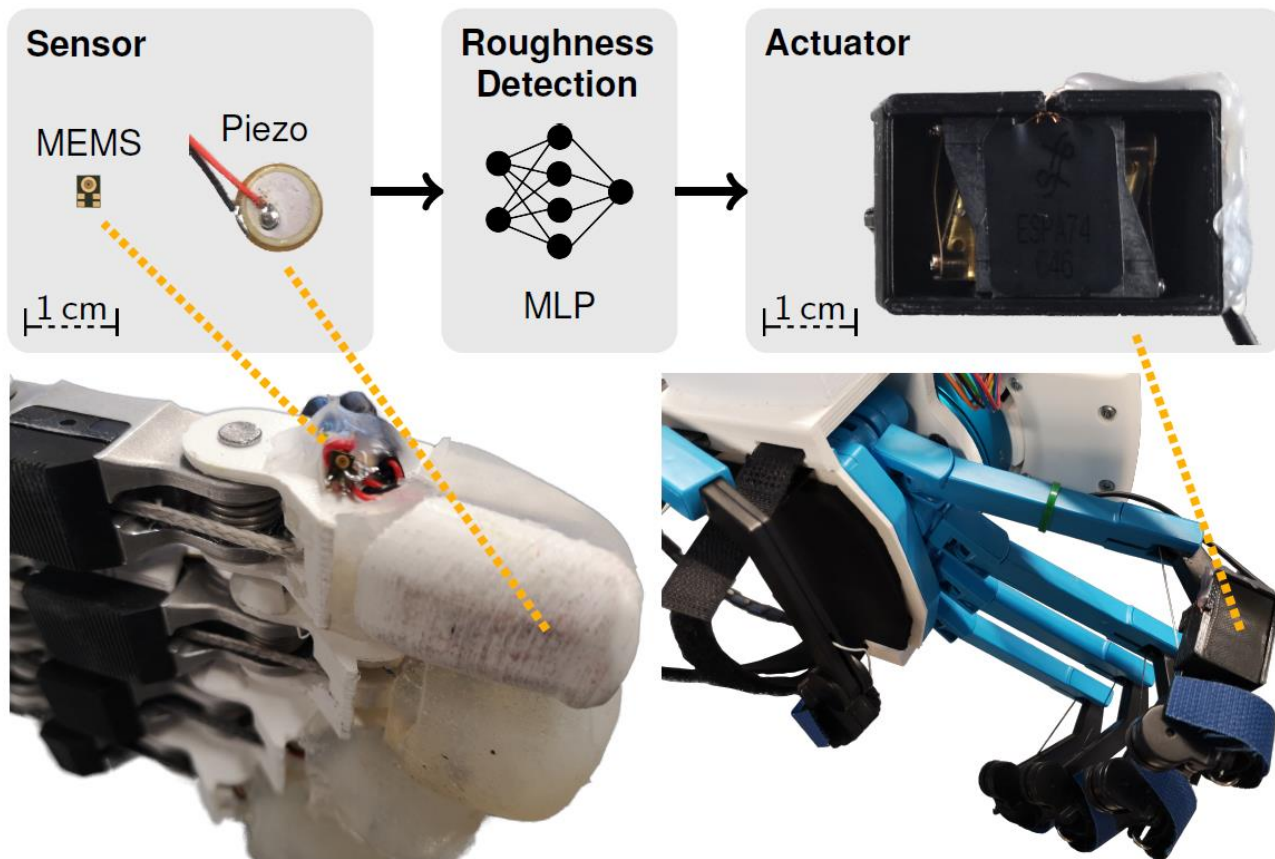


- Actuators on the hand exoskeleton



[Pätzold et al. arXiv:2303.07186]

# Roughness Perception



Data set of rough and smooth objects



[Pätzold et al. arXiv:2303.07186]



# Team NimbRo

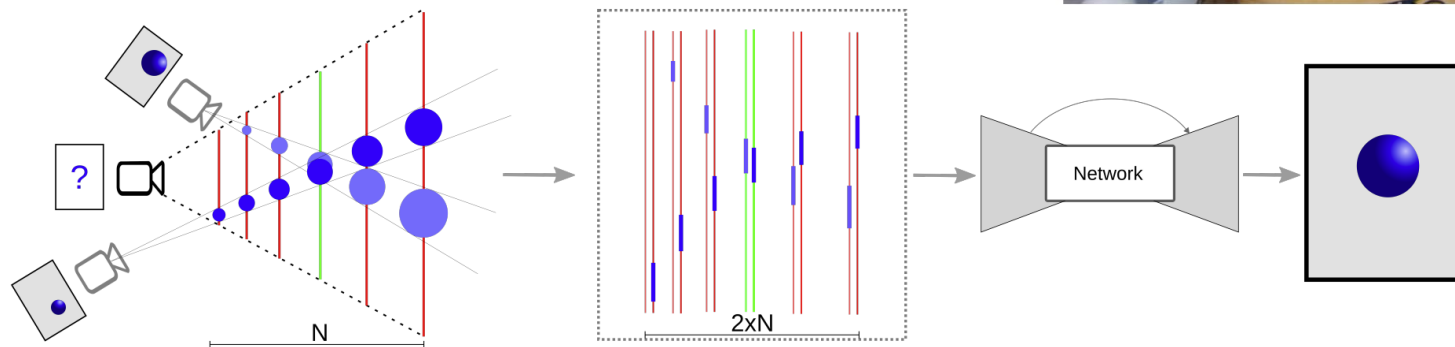
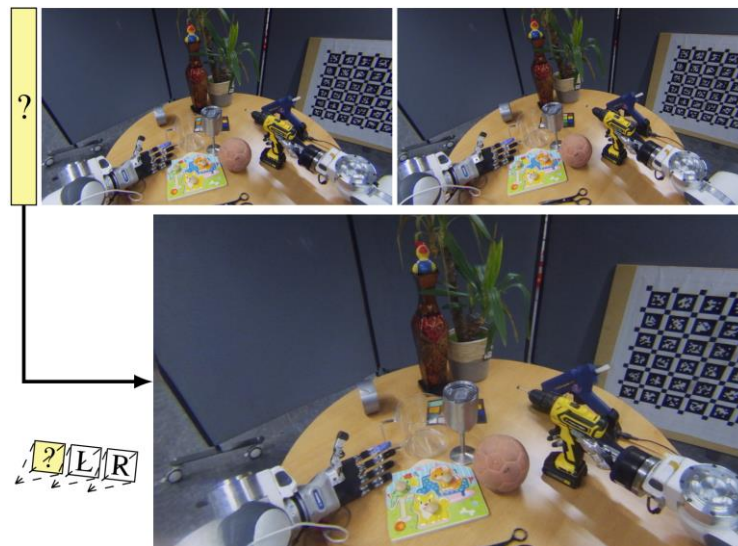


[Schwarz, Lenz et al.  
arXiv:2303.03297]



# FaDIV-Syn: Fast Depth-Independent View Synthesis

- Two input views
- Generate novel view from different pose
- Does not require depth
- Handles occlusions, transparency, reflectance, moving objects, ...



[Rochow et al. RSS 2022]

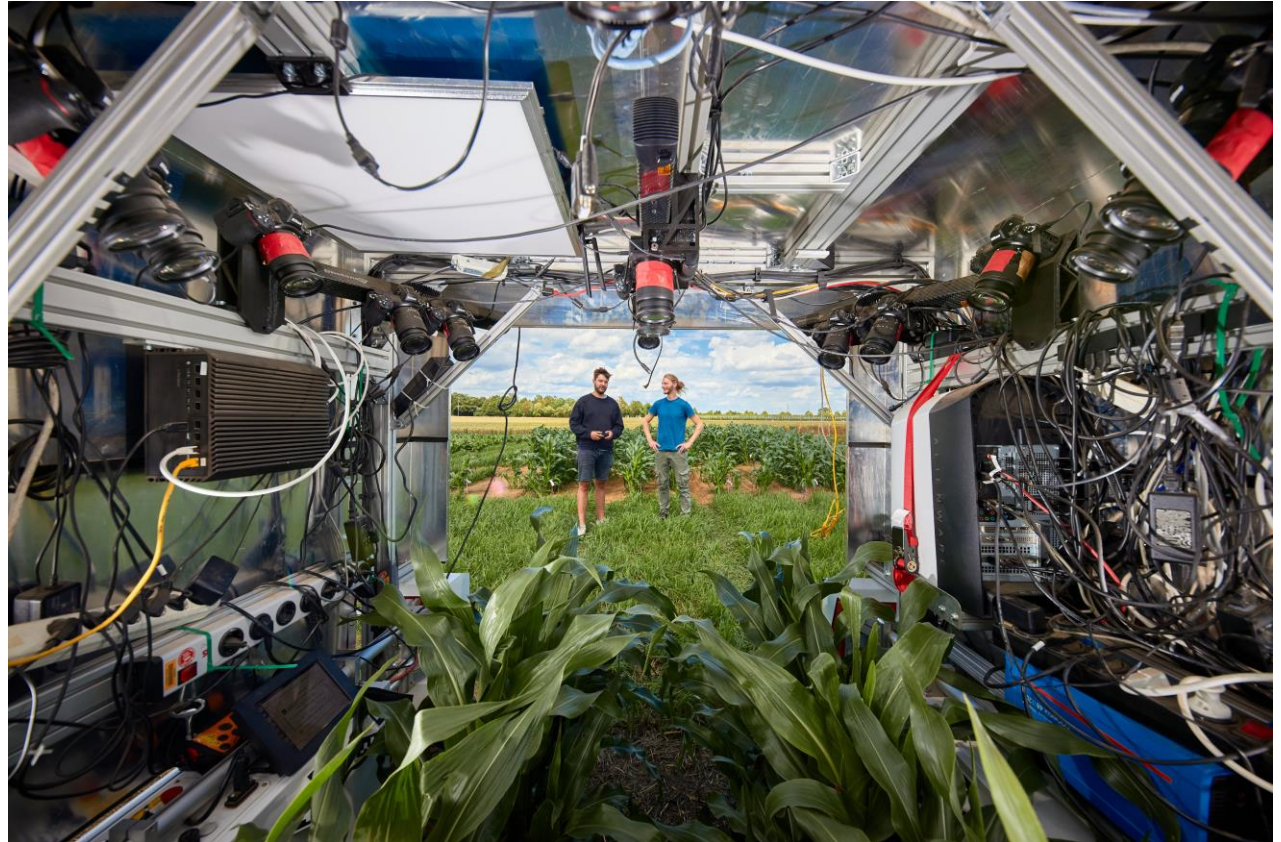
# FaDIV-Syn: Fast Depth-Independent View Synthesis

## Robot Teleoperation



# Multi-view Plant Reconstruction

- 14x Nikon Z7 DSLR camera
- 45 MP
- 64–25600 ISO
- 24-70 mm Lens



# Multi-view Plant Reconstruction

- Recovered camera poses and semi-dense point cloud through Multi-View-Stereo



# Implicit Surfaces on Permutohedral Lattices

- Geometry represented as Signed Distance Field (SDF)
- Color represented as a direction-dependent color field
- Transform SDF into radiance [1] and train similar to NeRF



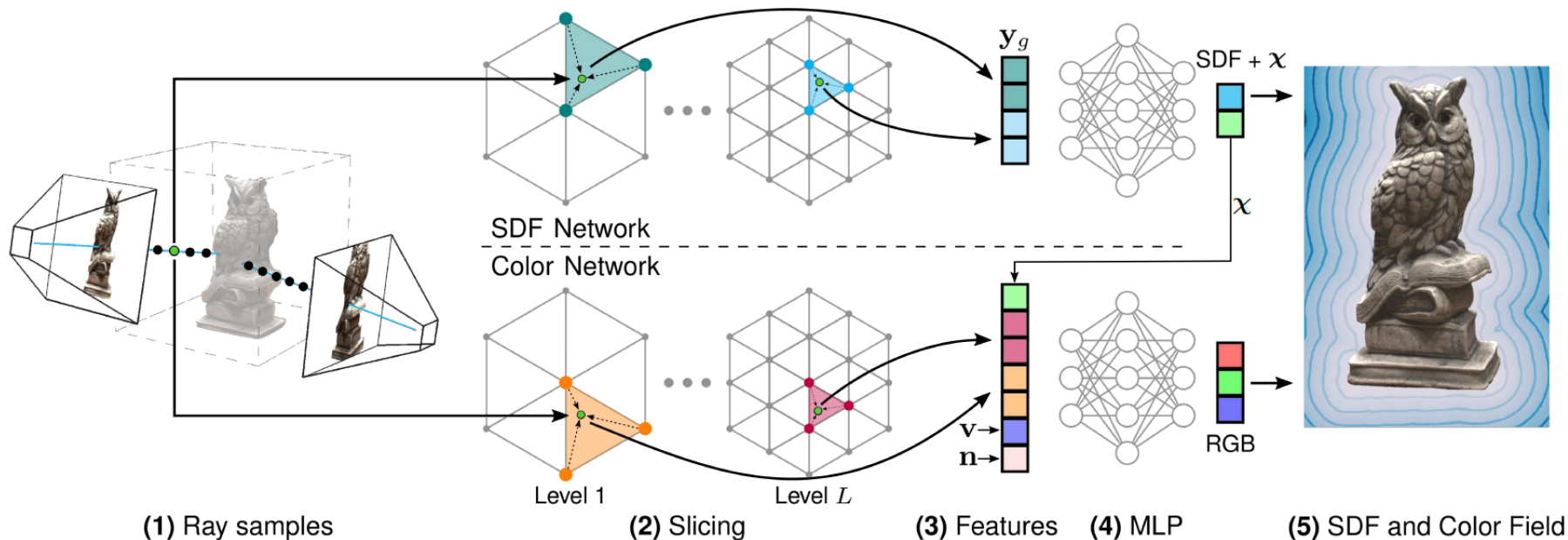
Geometry



Color at the zero level-set of the SDF

# Implicit Surfaces on Permutohedral Lattices

- Geometry represented as Signed Distance Field (SDF)
- Color represented as a direction-dependent color field
- Transform SDF into radiance [1] and train similar to NeRF



- InstantNGP with a Multiresolution Hash Encoding [2]
- Permutohedral lattice
- Small MLPs for SDF and color
- 25 M parameters
- 1 h training on Nvidia RTX 3090 GPU

[2] Müller et al. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding ACM Transactions on Graphics (SIGGRAPH 2022)

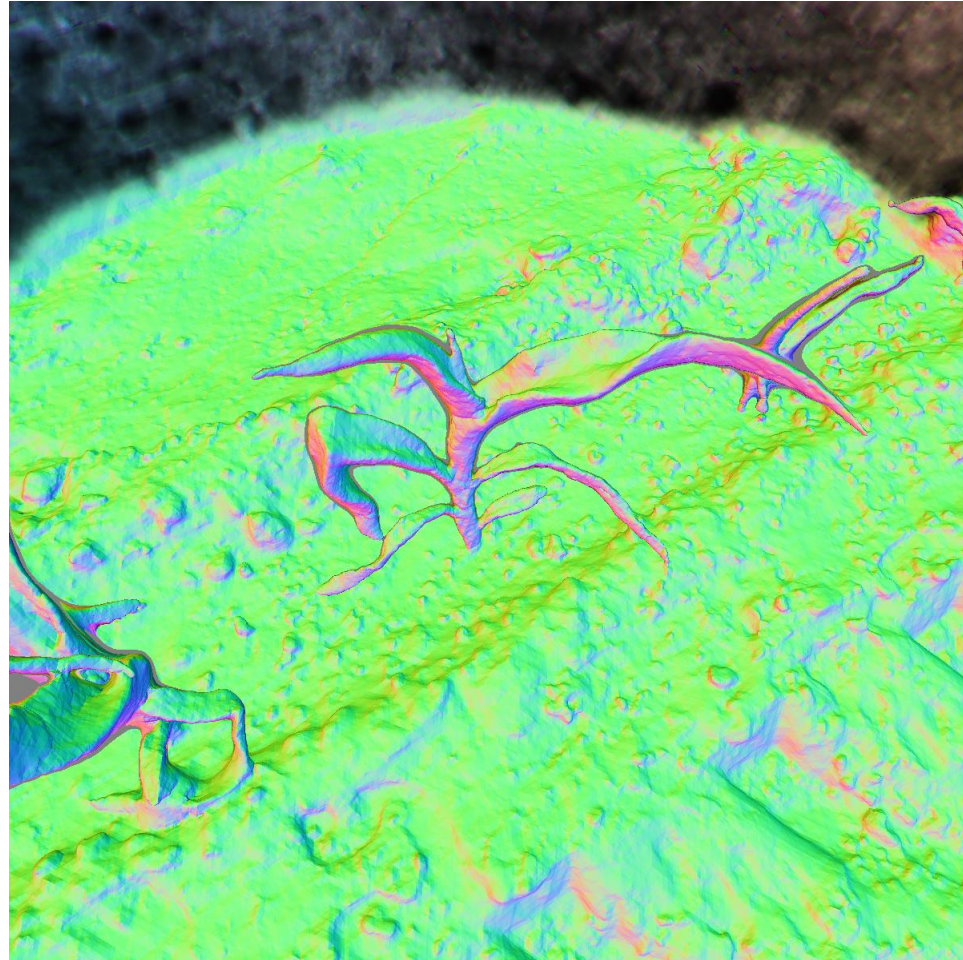
Surface normals



# Multi-view Plant Reconstruction

- InstantNGP with a Multiresolution Hash Encoding [2]
- Small MLPs for SDF and color
- 25 M parameters
- 1 h training on Nvidia RTX 3090 GPU

[2] Müller et al. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding ACM Transactions on Graphics (SIGGRAPH 2022)



Surface normals



# Multi-view Plant Reconstruction

- Rendered novel views



# Plant Reconstruction over Multiple Days



Volumetric renders through  
SDF + color

# Plant Reconstruction over Multiple Days



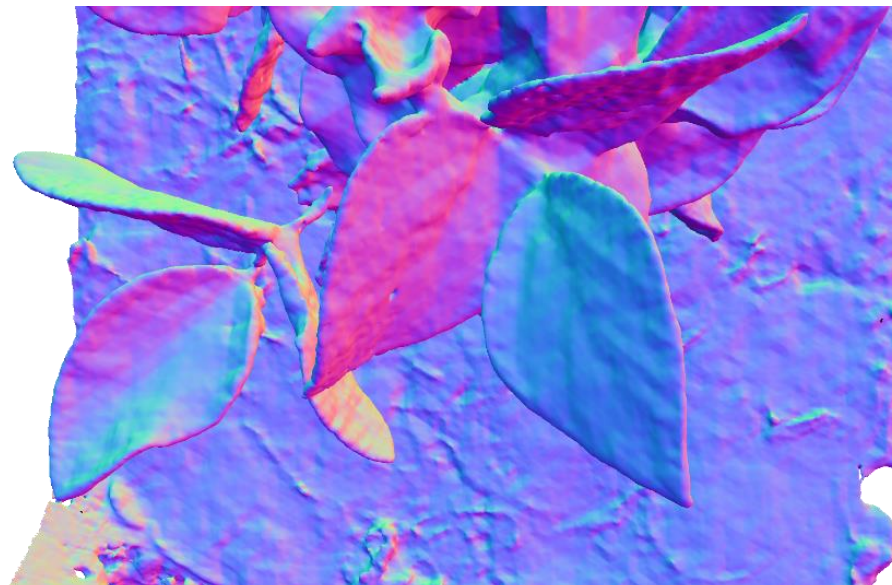
Predicted depth

# High Geometric and Texture Detail

- Marching cubes on the SDF to recover mesh
- Learnable texture to match color images
- Rendering in real time



Textured mesh



Mesh normal vector

# Conclusions

- Developed capable robotic systems for challenging scenarios
  - Bin picking
  - Humanoid soccer
  - Disaster response (UGV, UAV)
  - Plant reconstruction
- Challenges include
  - 4D semantic perception
  - High-dimensional motion planning
- Promising approaches
  - Prior knowledge (pretrained models, inductive bias)
  - Shared experience (fleet learning)
  - Shared autonomy (human-robot)
  - Instrumented environments

