

A white robotic arm is positioned in a laboratory setting. In the foreground, a black table holds several objects: a white cup, a blue cylinder, a white rectangular block, a yellow block with a red top, and a stack of three colored blocks (red, green, yellow). The background features a desk with multiple computer monitors, a whiteboard with handwritten notes and diagrams, and a window with blinds. A purple decorative shape is visible in the bottom right corner.

# The coupling of perception and interaction

For object discovery and understanding

Jen Jen Chung, Francesco Milano

14 October 2024





RPPL Lab

2023



2022




Autonomous Systems Lab



# Why dense object instance-aware scene reconstruction?



# Exploring interactions in an object-level map

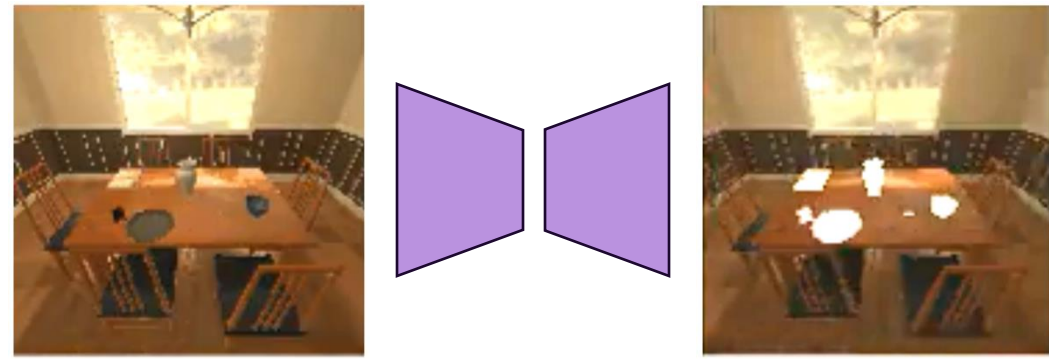


iTHOR simulation environment

- ✓ Robot pose
- ✓ RGB-D
- ✓ Instance mask
- ✓ Interaction result

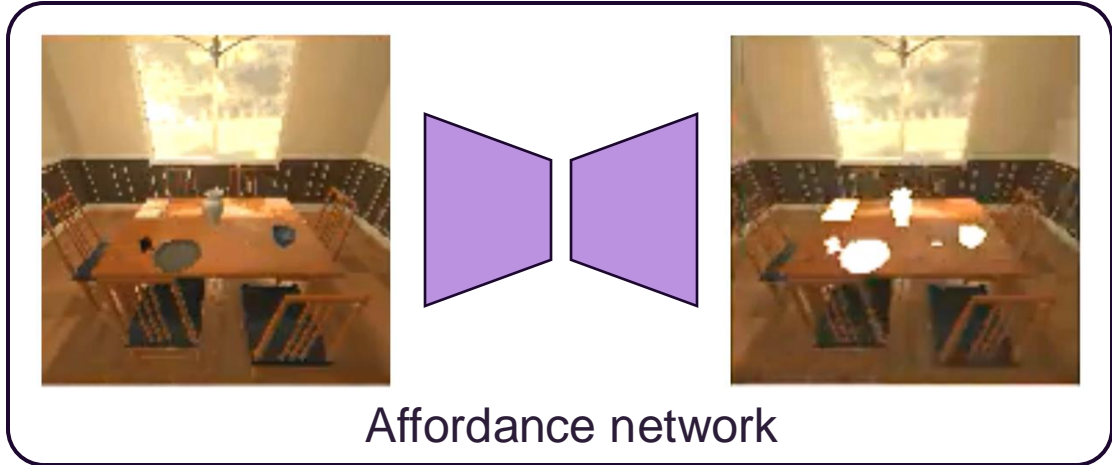
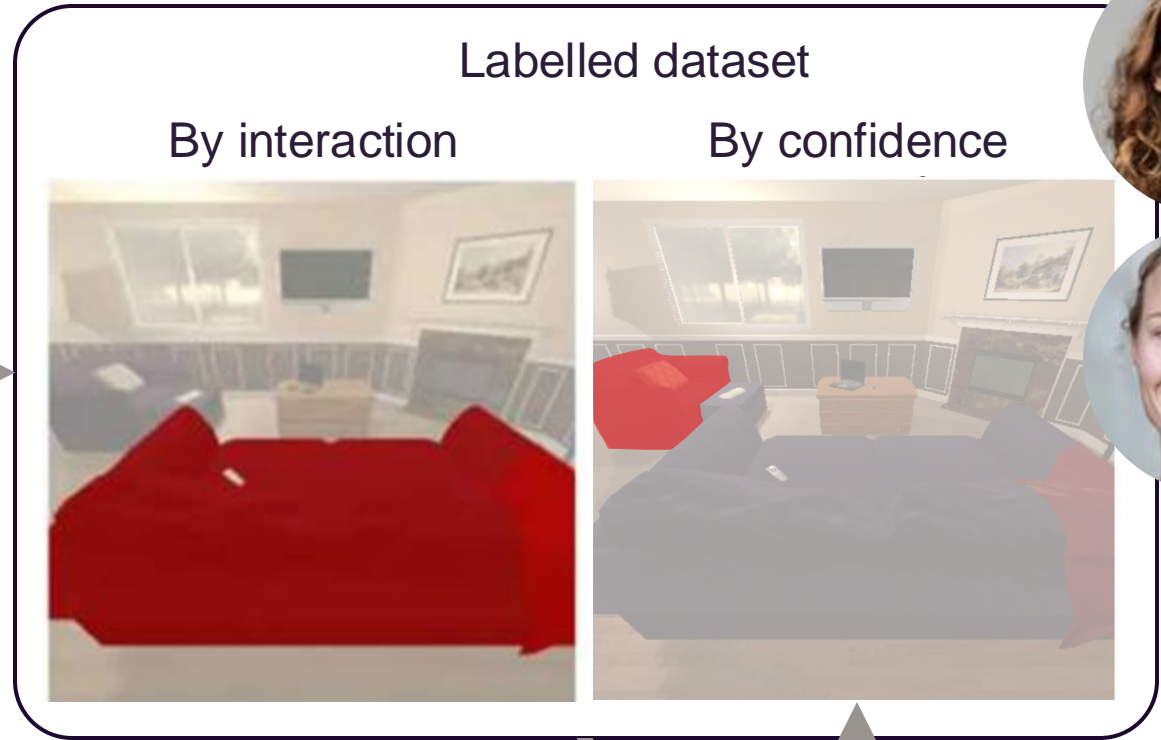


Object instance-level map



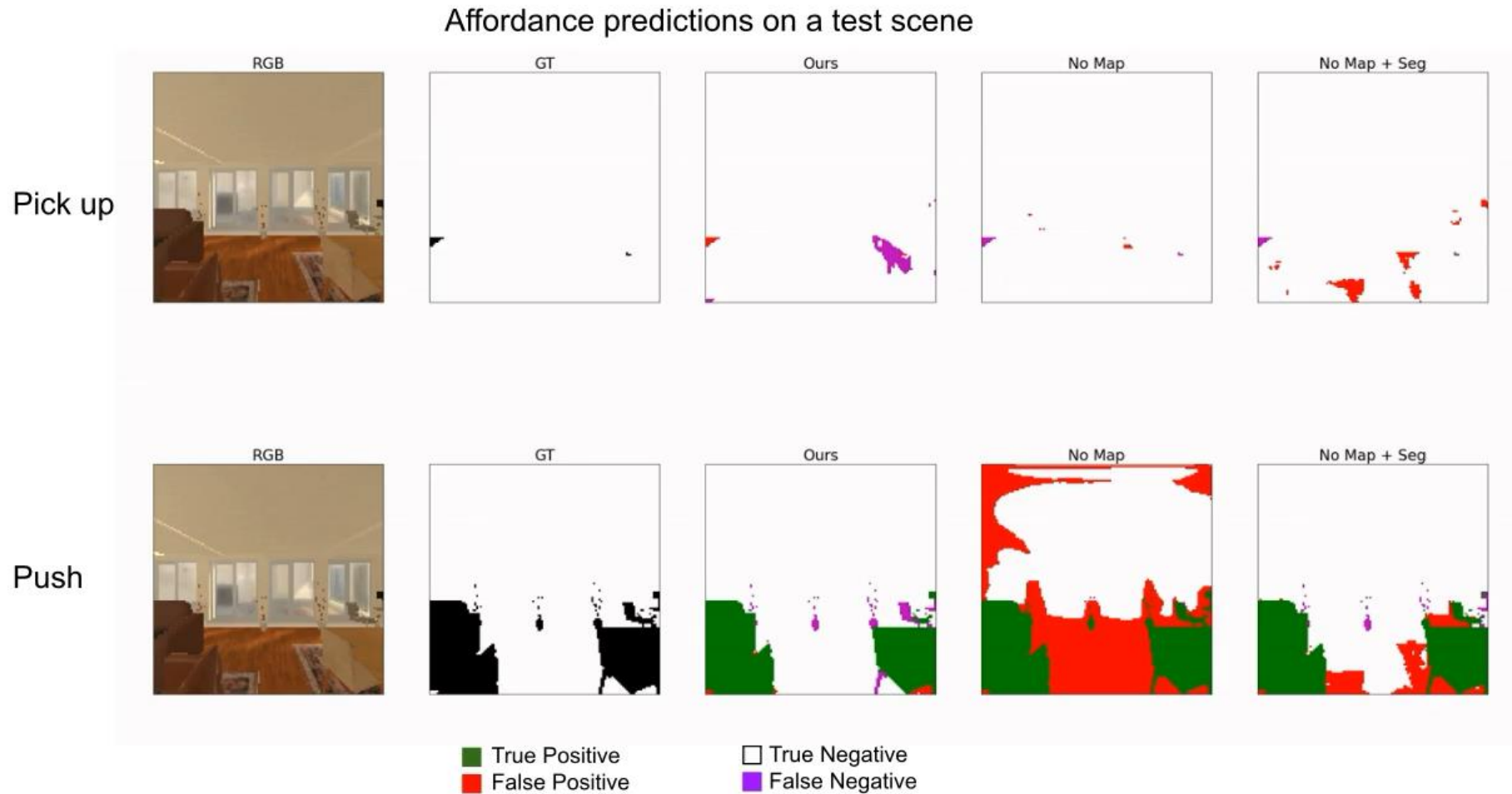
Affordance network

# Exploring interactions in an object-level map

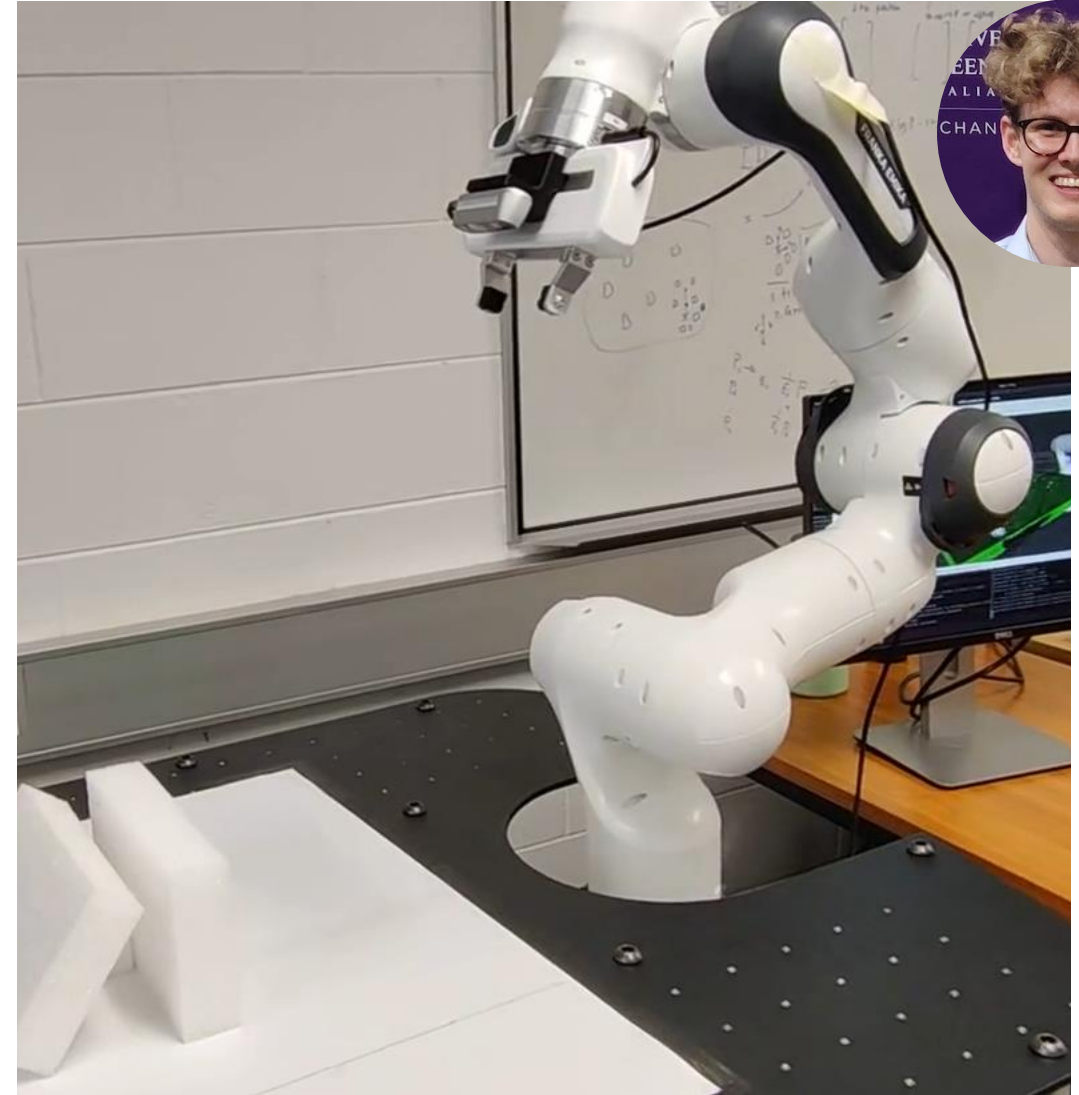
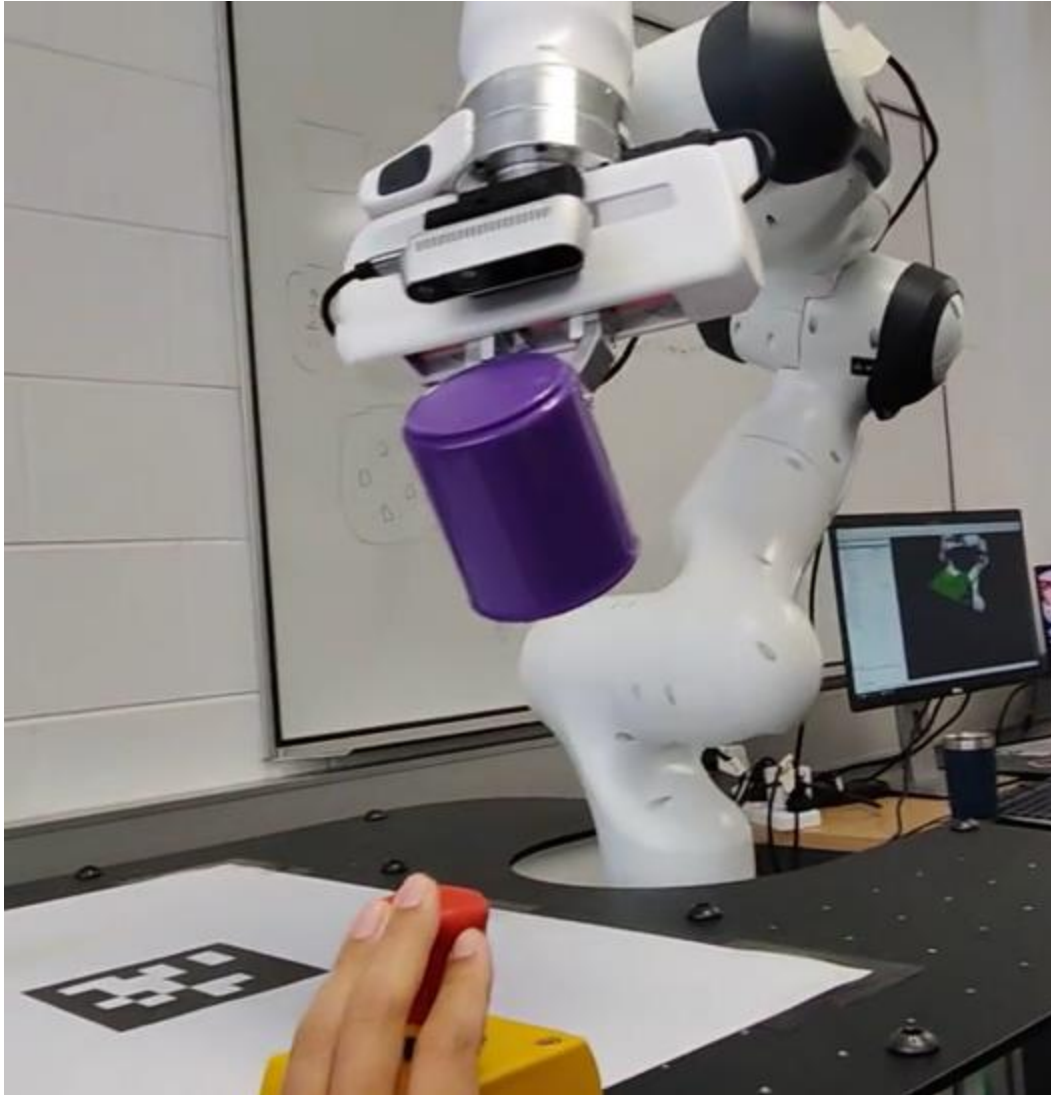




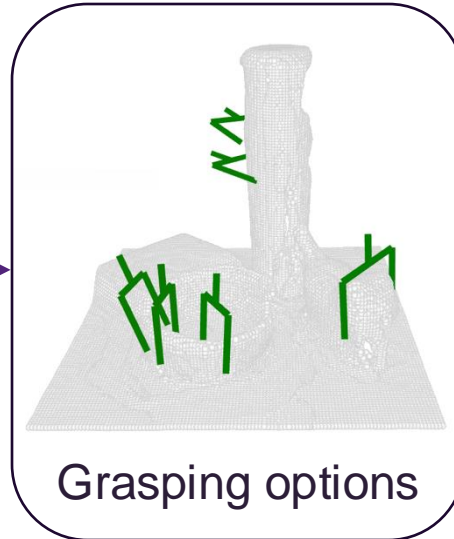
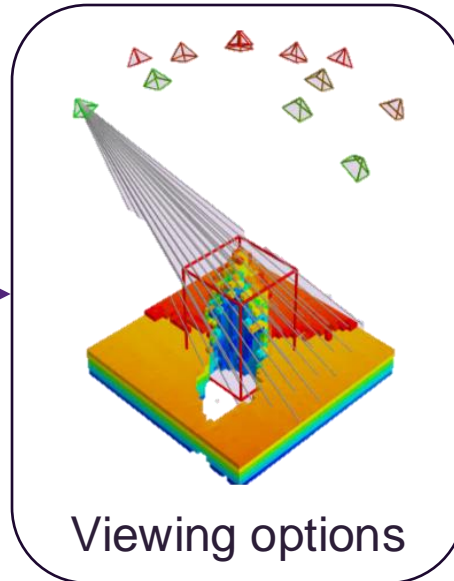
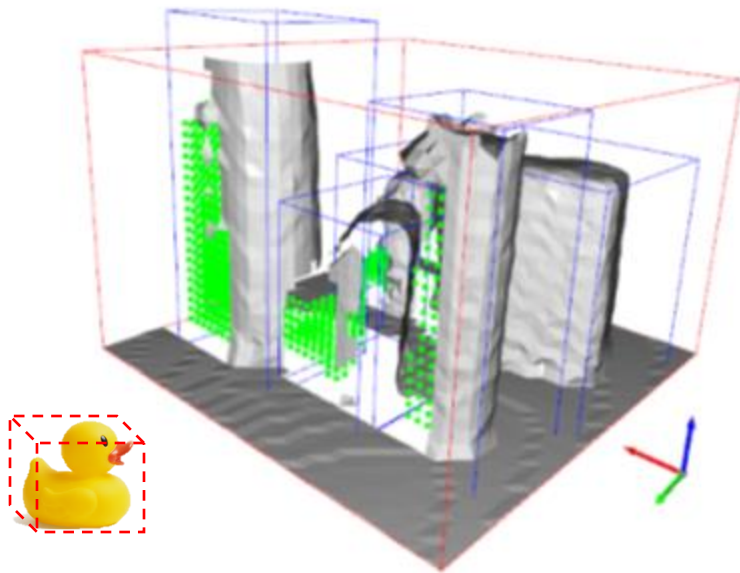
# Learned affordances from interactive exploration



# Finding and retrieving hidden objects



# To grasp or not to grasp?

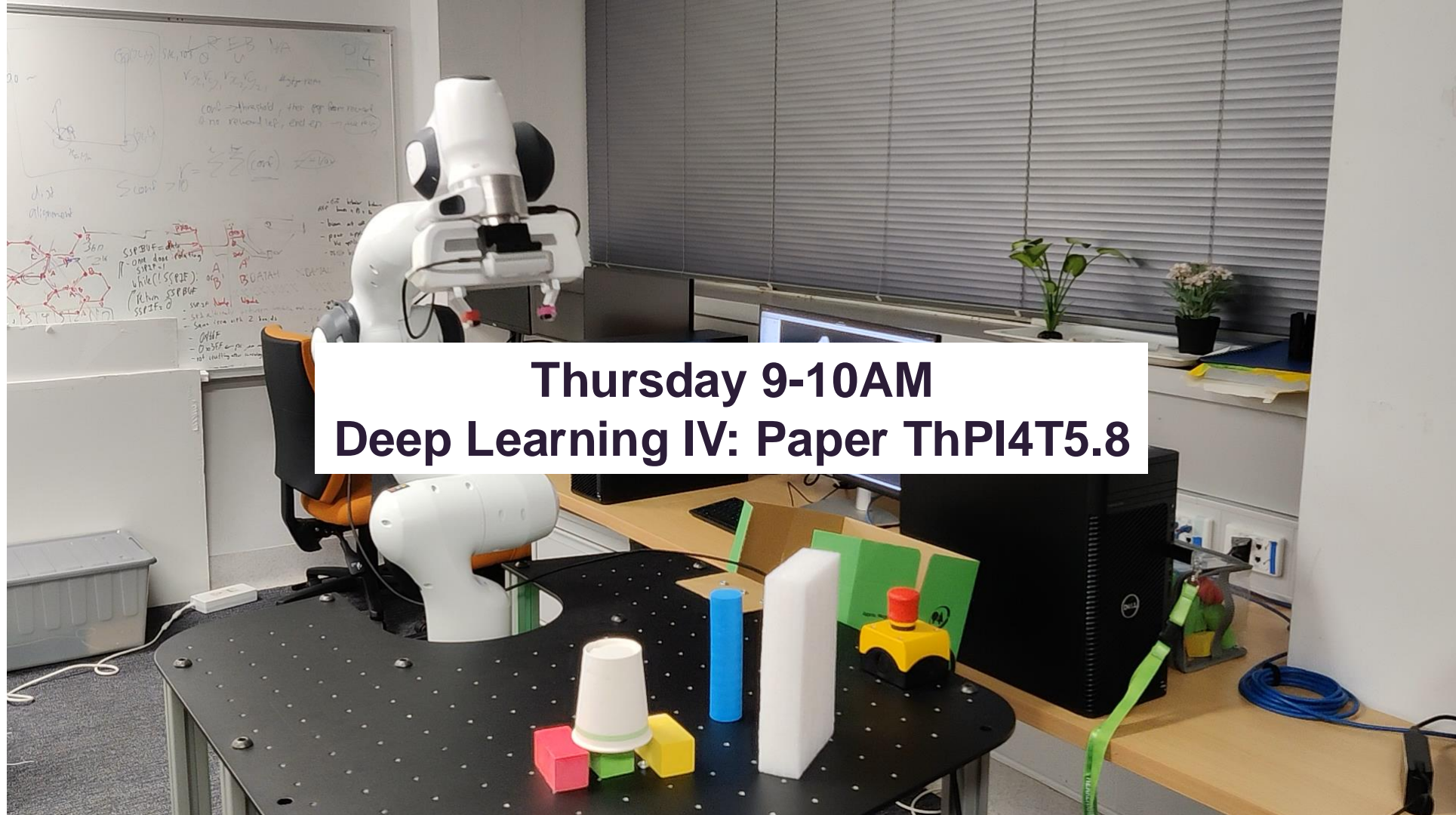


Action values learned via RL

$$R = \omega \frac{\Delta \text{ Possible target locations}}{\text{Possible target locations}} - t$$



# Active search and grasp in clutter



**Thursday 9-10AM**  
**Deep Learning IV: Paper ThPI4T5.8**



# **Object-level representations for robotic interaction**



# Harmony: Assistive robots for healthcare



# Harmony: Assistive robots for healthcare





# Harmony: Assistive robots for healthcare



Automation of hospital bioassay sample flow



**Harmony**  
Assistive robots for healthcare

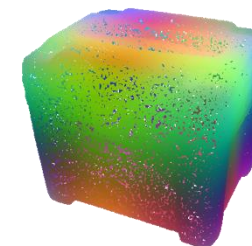
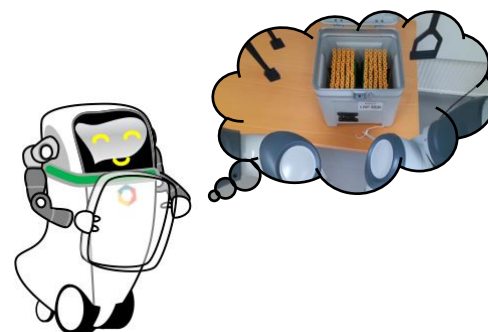


# Perception in support of robotic interaction

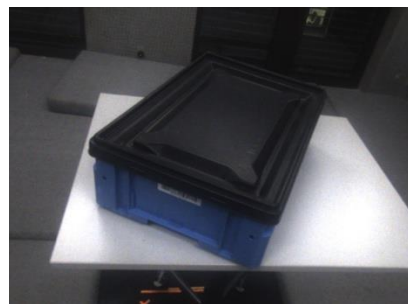
Desiderata:



**Accurate reconstruction** (geometry, appearance)



**Flexibility: Encode task-specific properties**



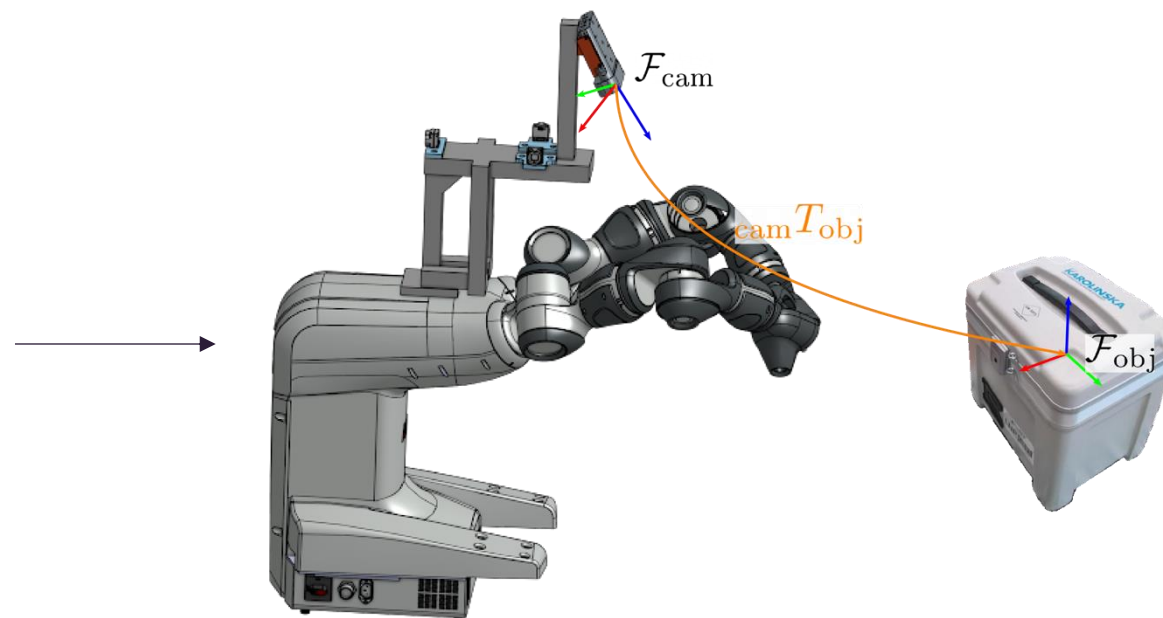
**Ability to easily incorporate new representations**

How? Our hypothesis: **Neural Fields + Neural Rendering**



# Perception in support of robotic interaction

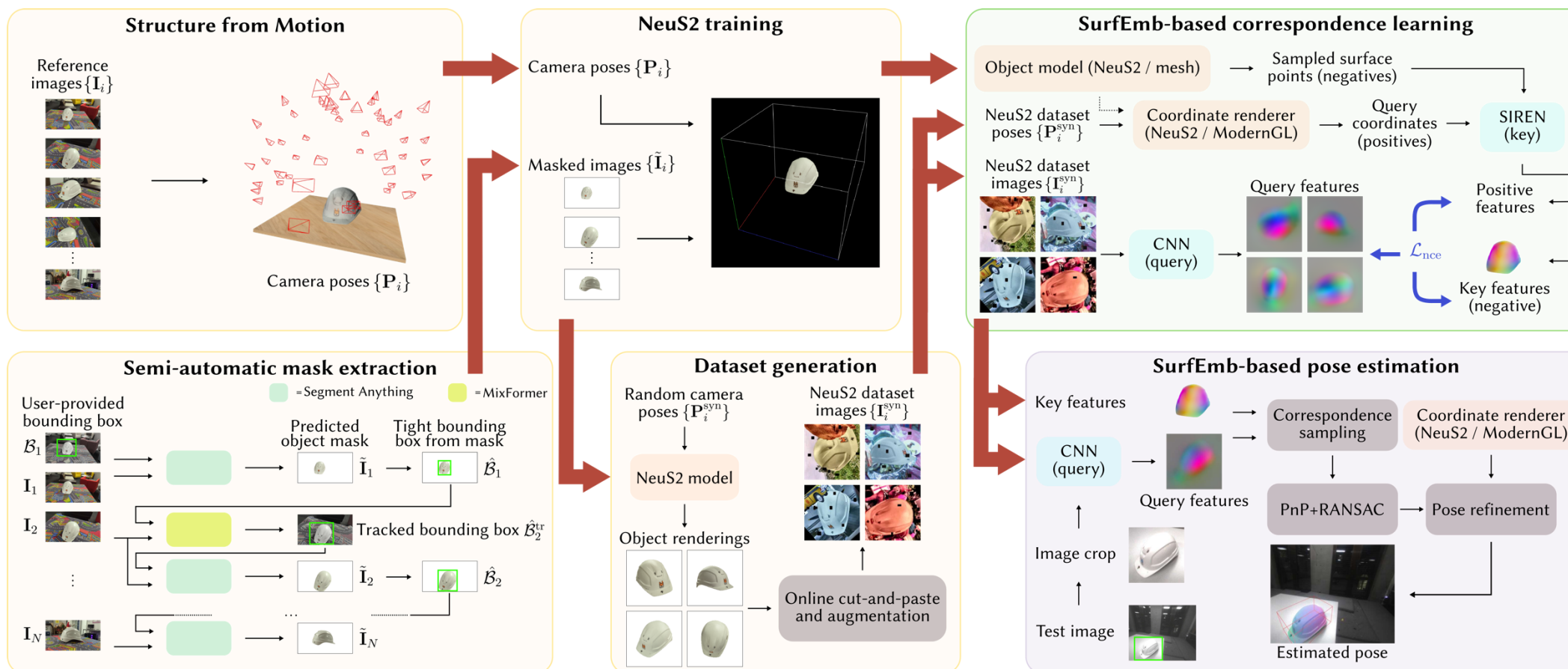
One example task: 6-DoF Object Pose Estimation



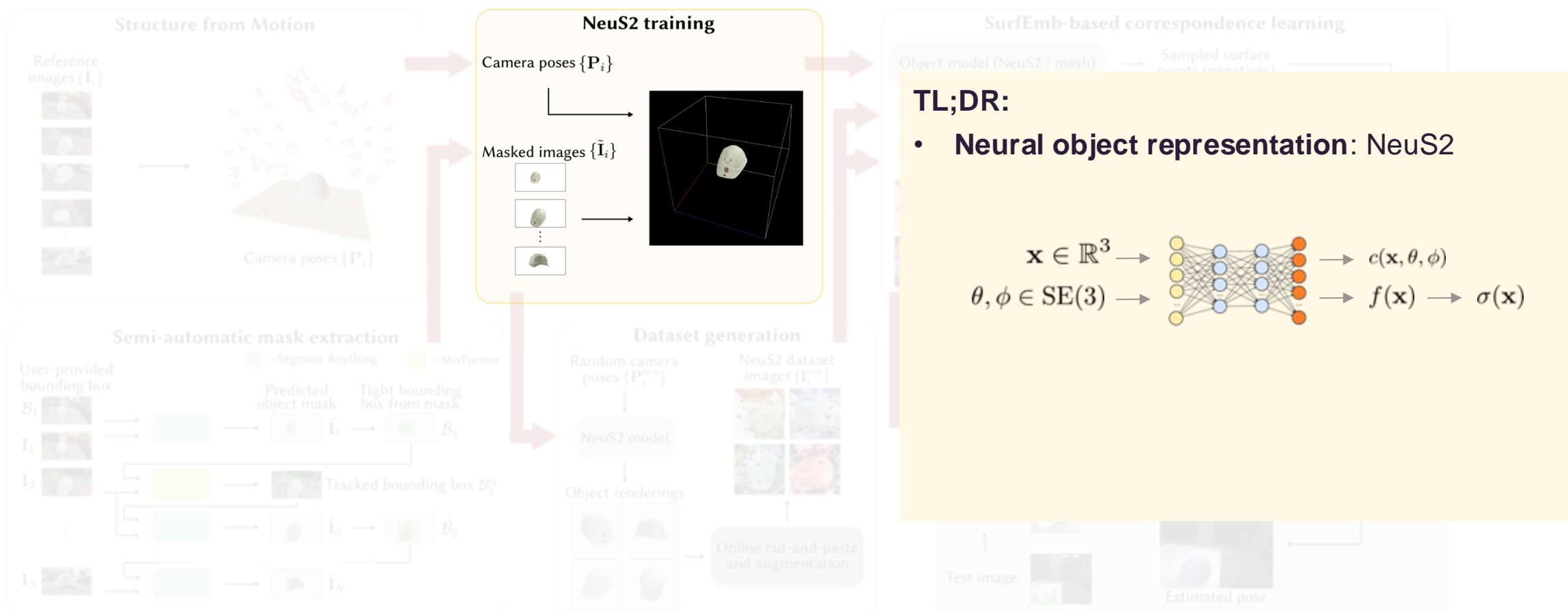
■ State-of-the-art approaches rely on **textured CAD models** and **photorealistic synthetic datasets** (PBR)

How can neural fields and neural rendering help?

# Neural fields for object pose estimation – NeuSurfEmb

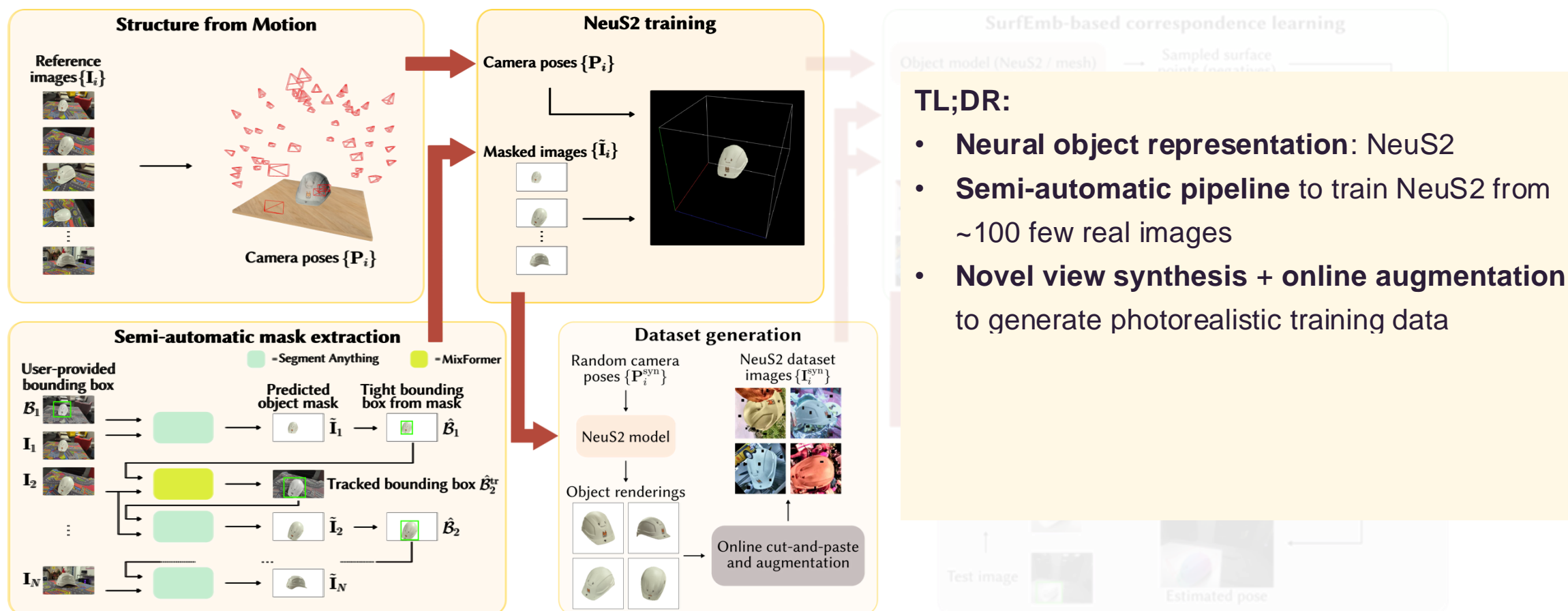


# Neural fields for object pose estimation – NeuSurfEmb

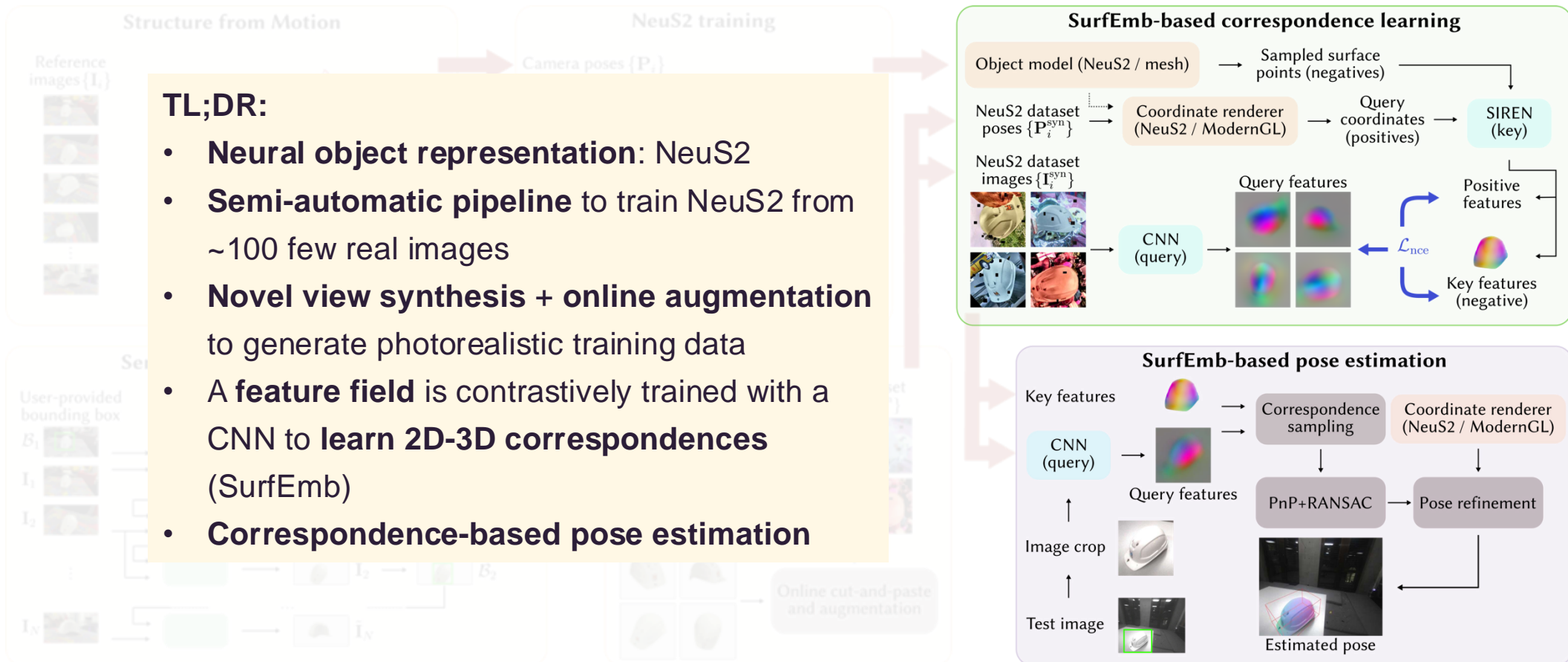




# Neural fields for object pose estimation – NeuSurfEmb



# Neural fields for object pose estimation – NeuSurfEmb





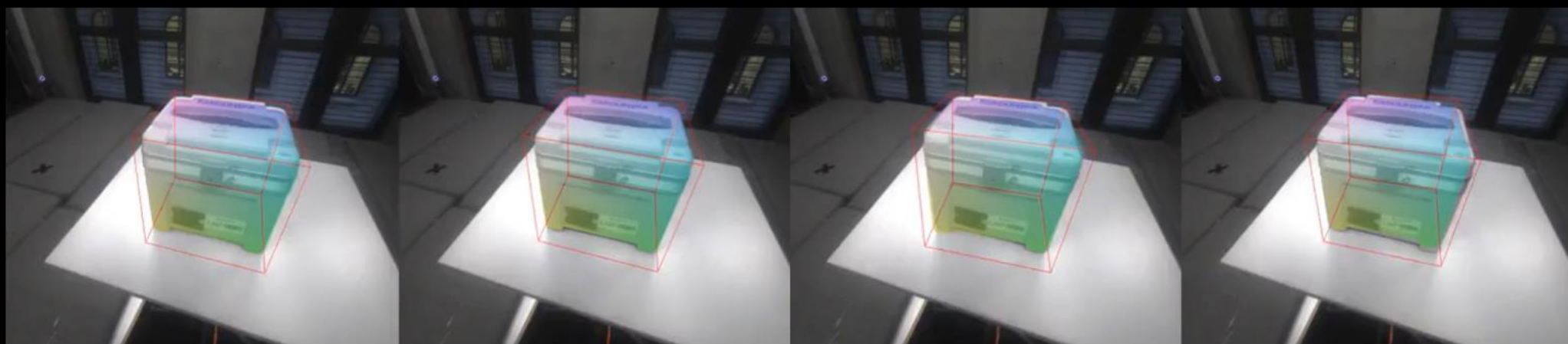


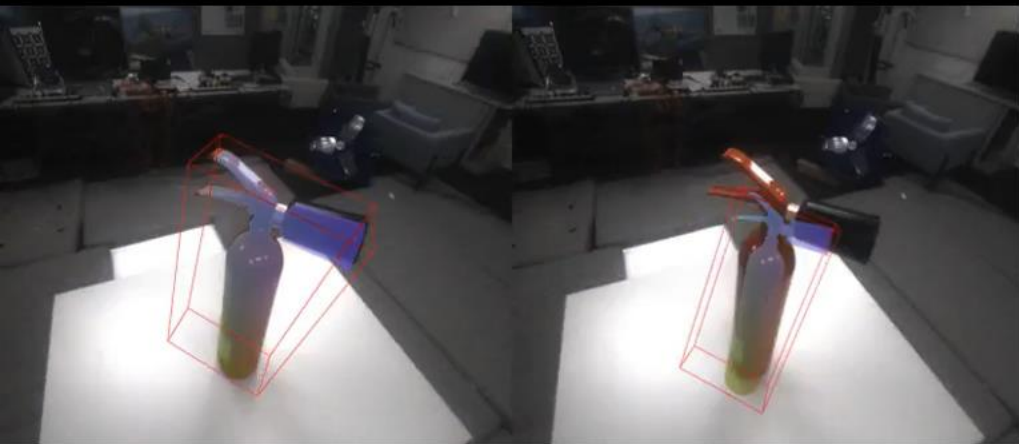
NeuSurfEmb

OnePose++

Gen6D (with tracking)

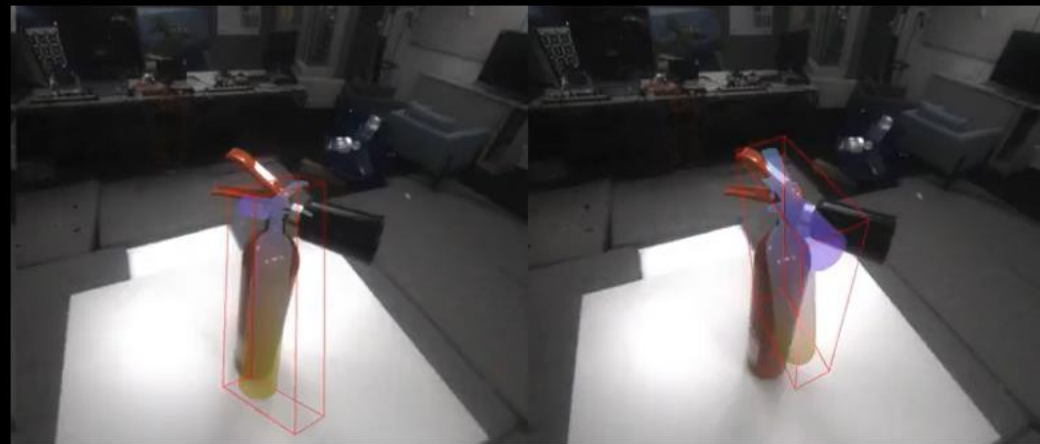
Gen6D (w/o tracking)





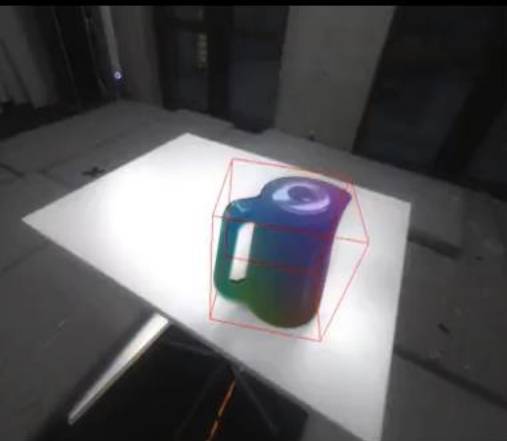
NeuSurfEmb

OnePose++ (w/o  
tracking, orig. recrop.)

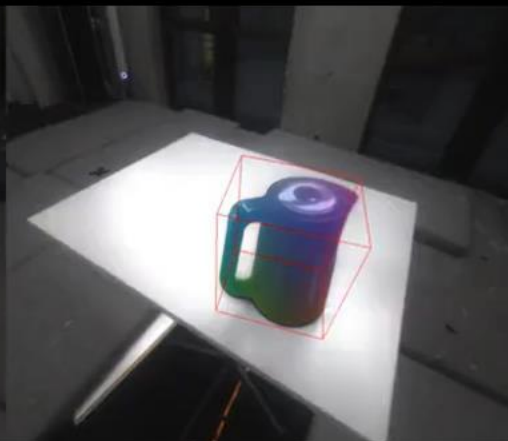


Gen6D (with tracking)

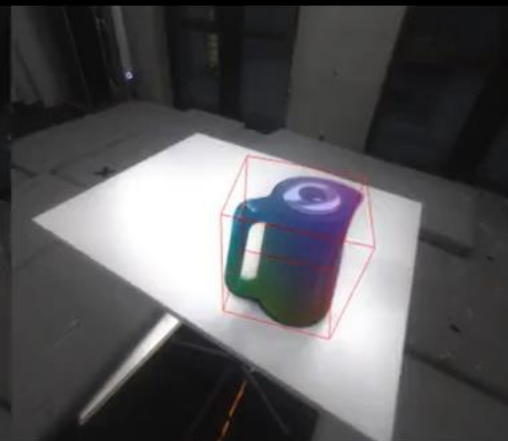
Gen6D (w/o tracking)



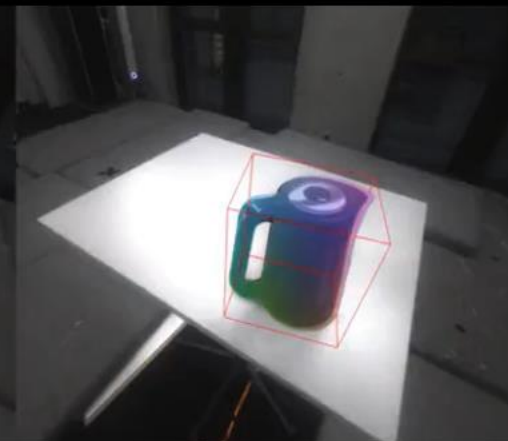
NeuSurfEmb



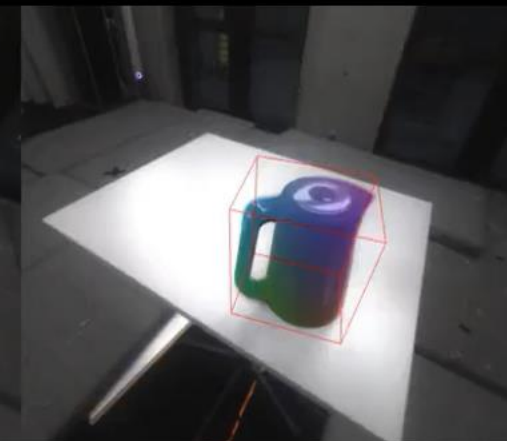
OnePose++ (w/o  
tracking, orig. recrop.)



OnePose++ (w/o  
tracking, prop. recrop.)

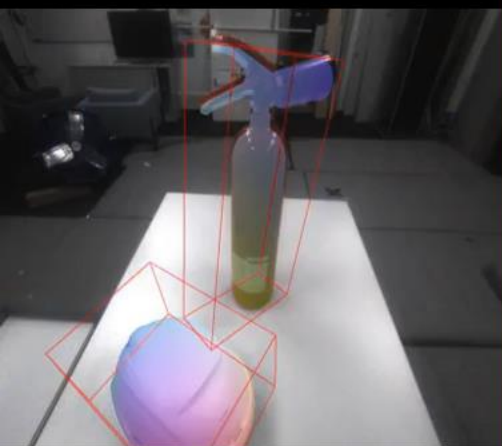


Gen6D (with tracking)

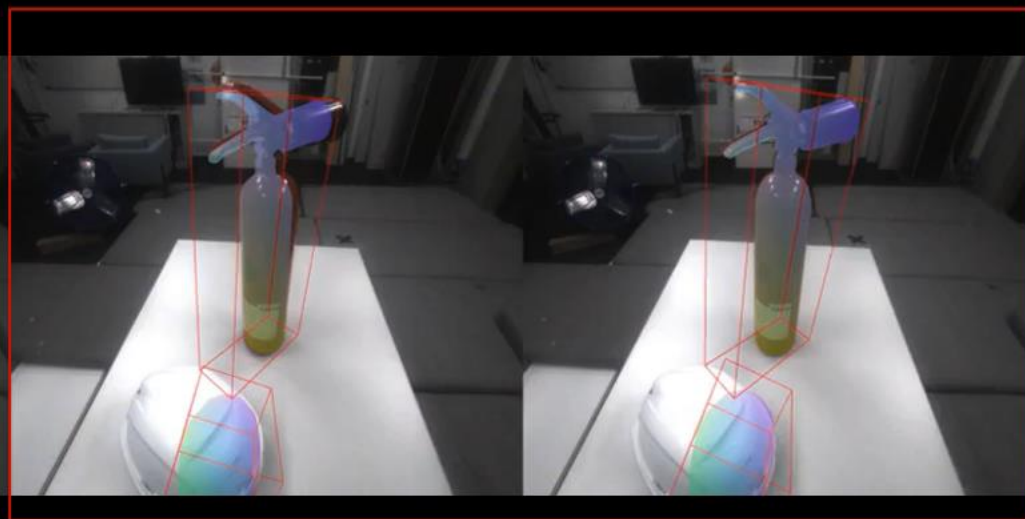


Gen6D (w/o tracking)

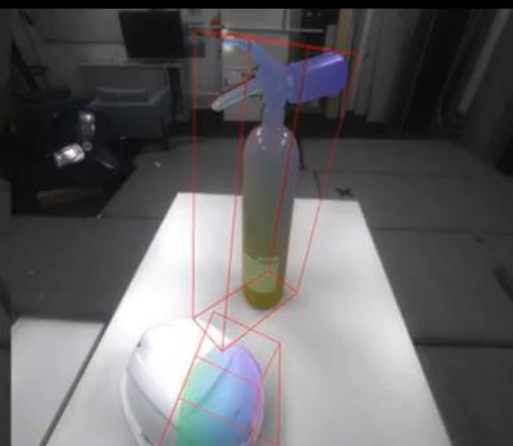




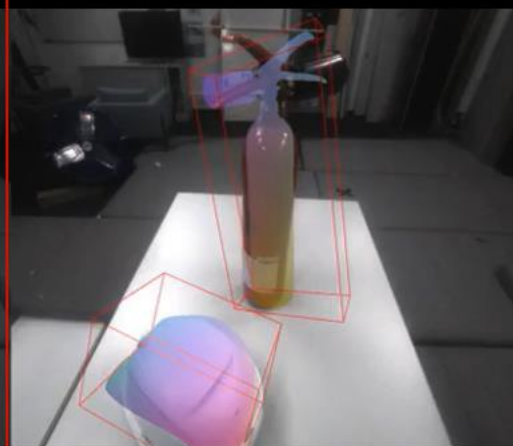
NeuSurfEmb



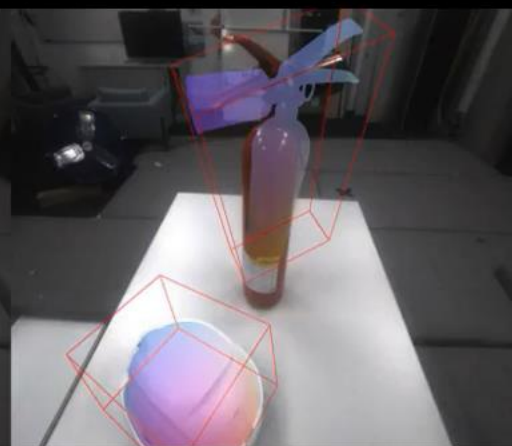
OnePose++ (w/o  
tracking, orig. recrop.)



OnePose++ (w/o  
tracking, prop. recrop.)

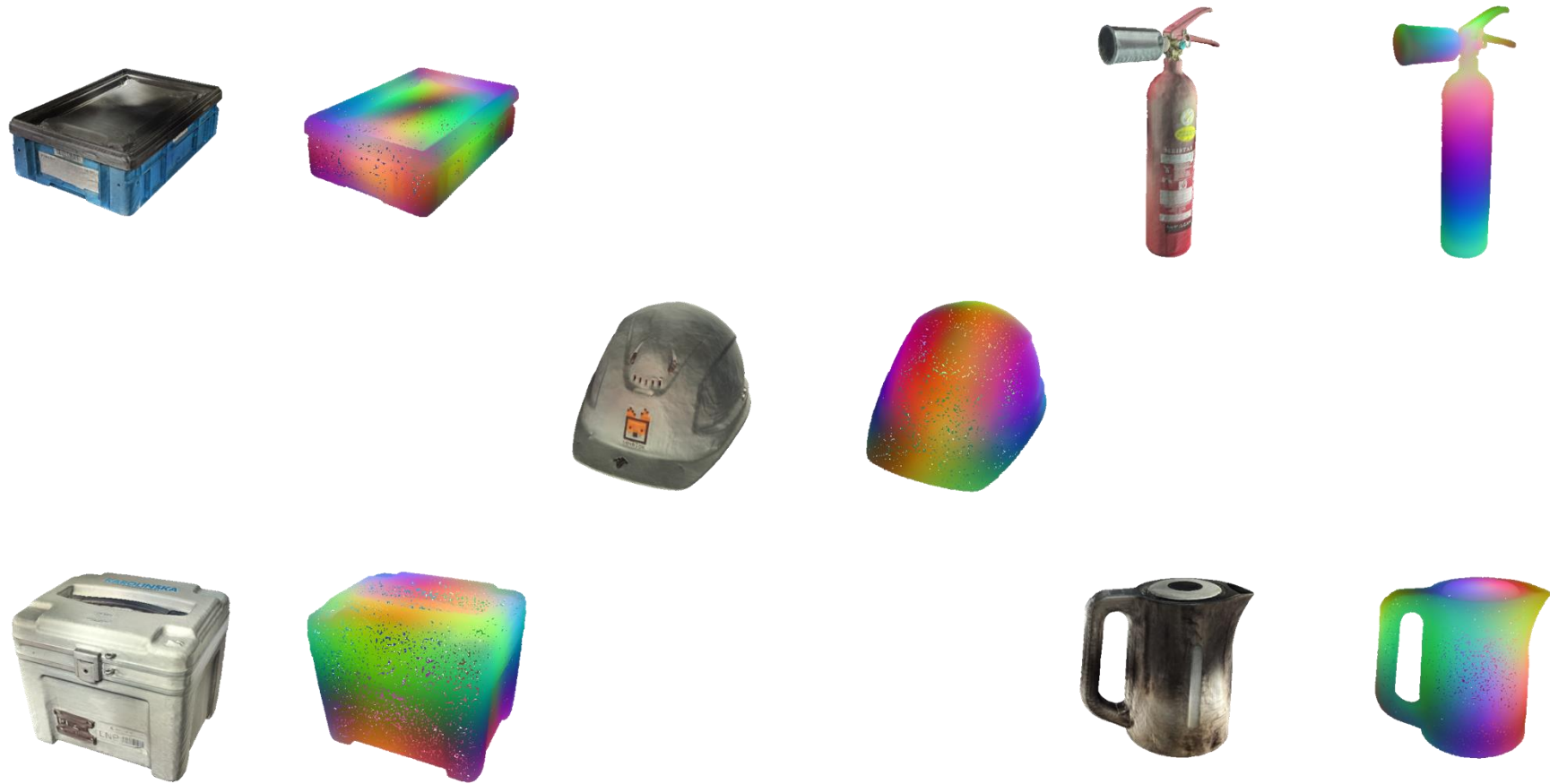


Gen6D (with tracking)



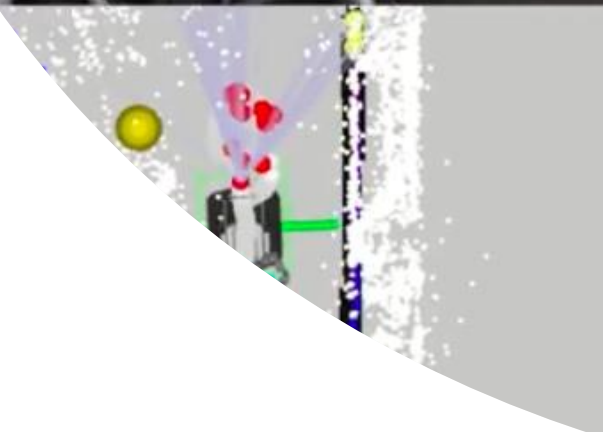
Gen6D (w/o tracking)

The low amount of texture generally causes less accurate predictions for OnePose++



## What is the future of object representations for robotics?

- How can we form object representations even more efficiently?
- What type of properties should we additionally incorporate?
- Is an explicit database needed or will implicit, large-scale priors be the future?
- ...



## Jen Jen Chung

Robotic Perception, Planning and Learning Lab, UQ  
jenjen.chung@uq.edu.au

## Francesco Milano

Autonomous Systems Lab, ETHZ  
francesco.milano@mavt.ethz.ch