# Perceive, Plan, Act and Learn: Towards Intelligent Robots in Human Environments
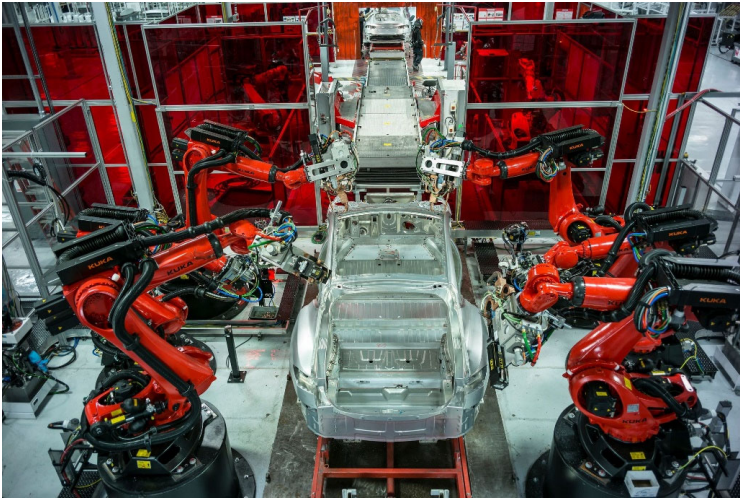
Yu Xiang

Senior Research Scientist

NVIDIA Research

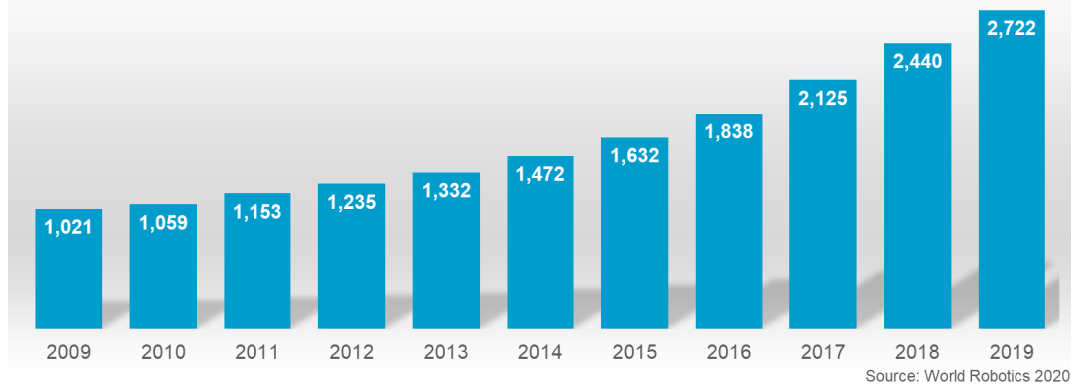# Robots in Factories and Warehouses



Welding and Assembling



Material Handling



Delivering



**Operational stock of industrial robots - World**
1,000 units

| Year | Units |
|------|-------|
| 2009 | 1,021 |
| 2010 | 1,059 |
| 2011 | 1,153 |
| 2012 | 1,235 |
| 2013 | 1,332 |
| 2014 | 1,472 |
| 2015 | 1,632 |
| 2016 | 1,838 |
| 2017 | 2,125 |
| 2018 | 2,440 |
| 2019 | 2,722 |

Source: World Robotics 2020

# Current Robots in Human Environments


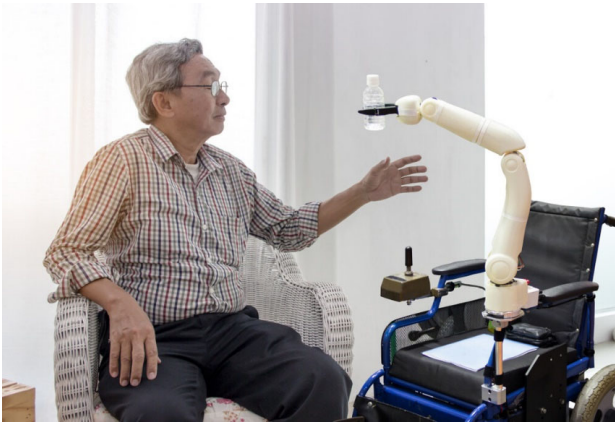Cleaning Robots


Telepresence Robots


Smart Speakers

How can we have more powerful robots assisting people at homes or offices?
- Mobile manipulators
- Humanoids



3

# Future Intelligent Robots in Human Environments


Senior Care


Assisting


Serving


Cooking


Cleaning


Dish washing

4

# Why Bringing Robots to Human Environments is Challenging?

Closed World: Factories & Warehouses



Open World: Human environments



- Structured environments
- Single tasks

- Unstructured and dynamic environments
- Various tasks

# Why Bringing Robots to Human Environments is Challenging?
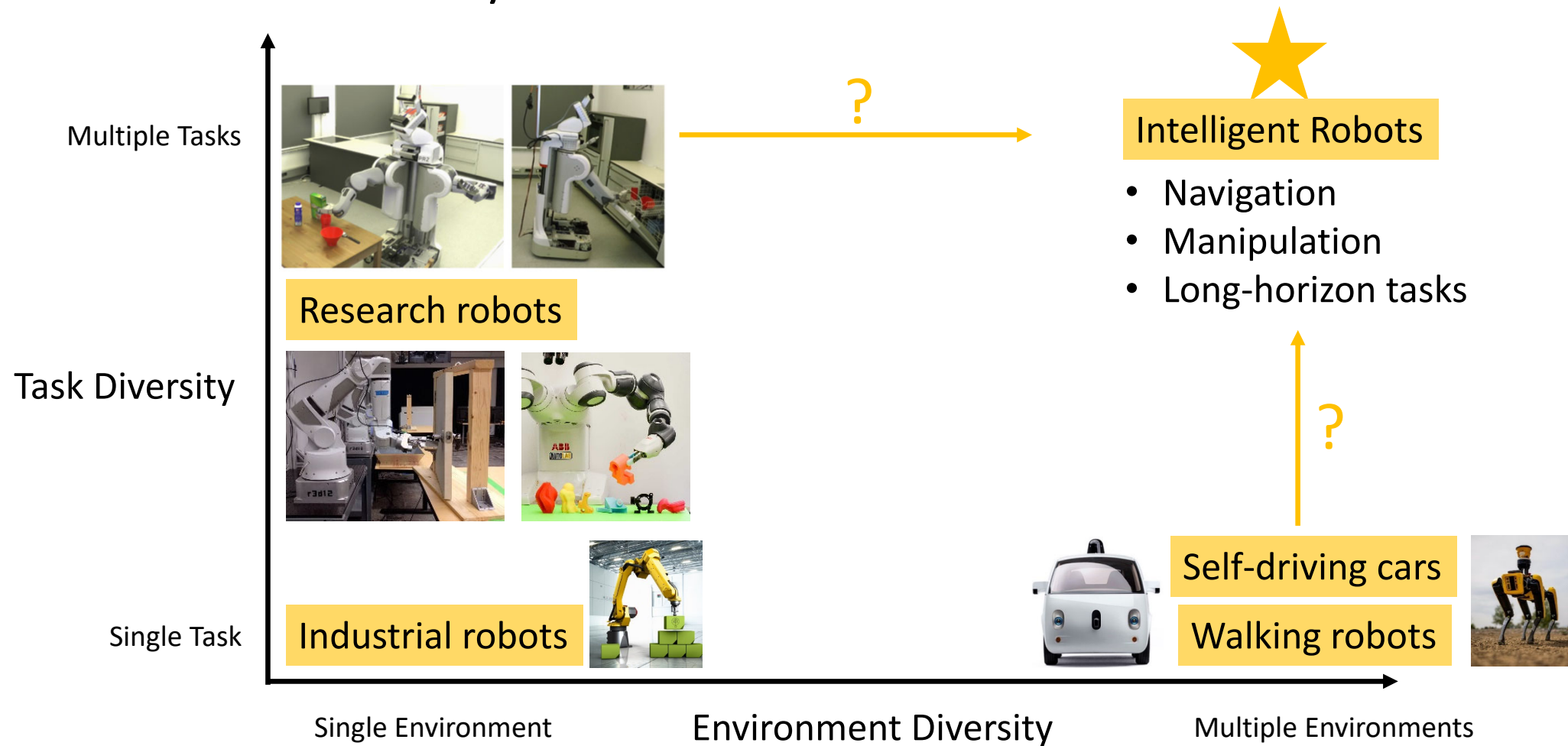
Example: Picking up a mug
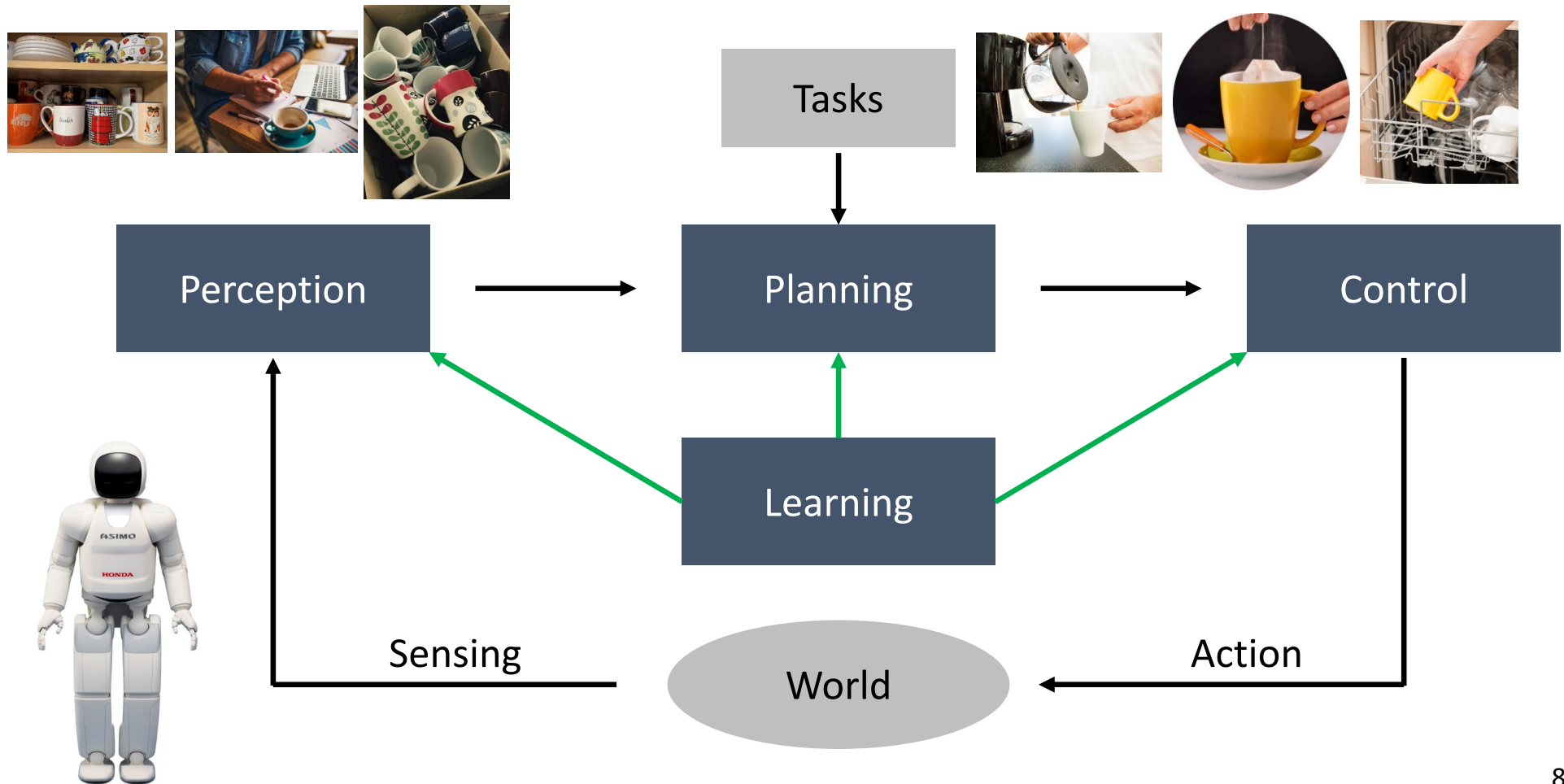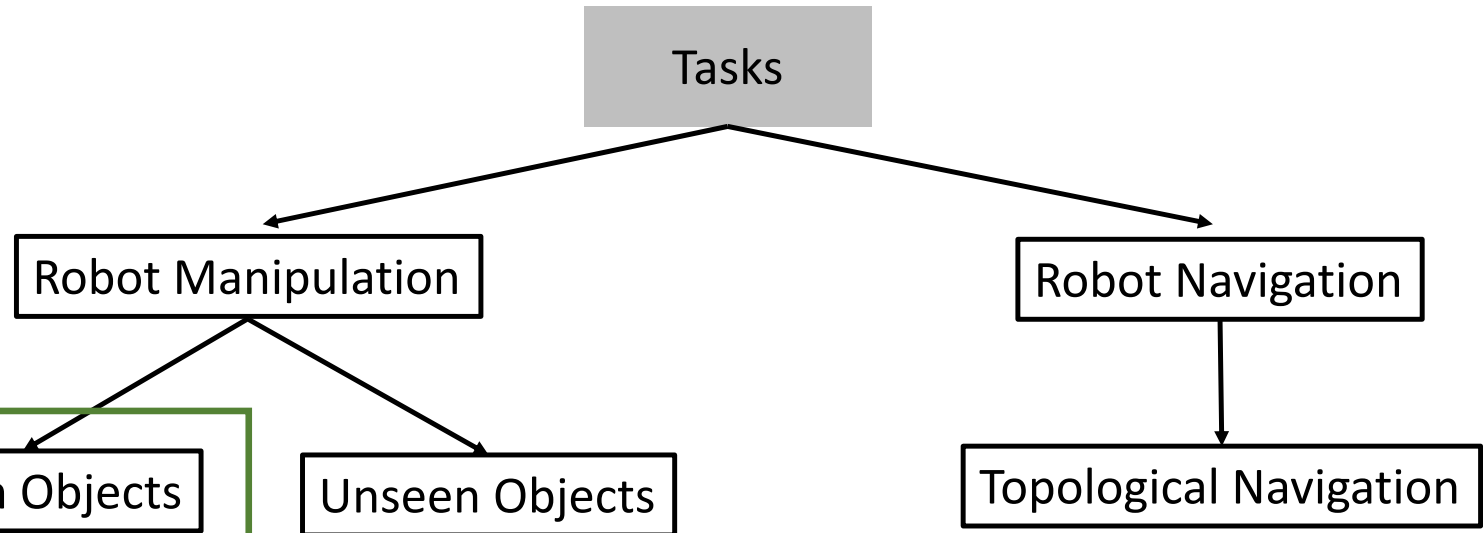

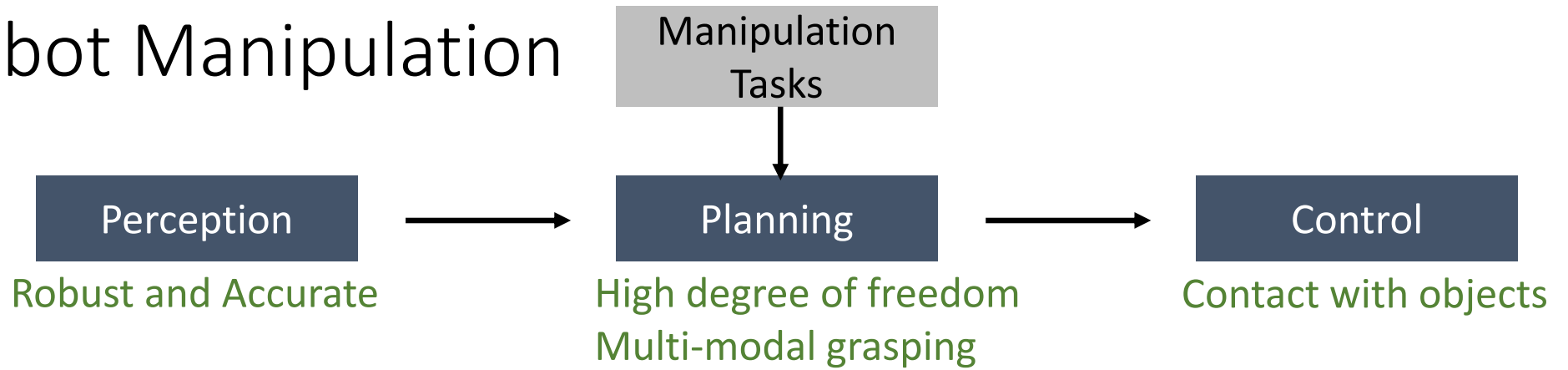
4X

Our Lab

**Environment Diversity**



**Task Diversity**



6

# Robot Autonomy



Task Diversity (vertical axis: Single Task → Multiple Tasks)

Environment Diversity (horizontal axis: Single Environment → Multiple Environments)

Research robots

Industrial robots

Self-driving cars

Walking robots

Intelligent Robots
- Navigation
- Manipulation
- Long-horizon tasks

?

?

# The Perception, Planning and Control Loop

# Outline



Tasks

Robot Manipulation

Robot Navigation

Known Objects

Unseen Objects

Topological Navigation

# Robot Manipulation

Manipulation Tasks

Perception → Planning → Control

Robust and Accurate

High degree of freedom
Multi-modal grasping

Contact with objects

Sensed image

Planning scene

Real world execution

2X

# Perception: Model-based 6D Object Pose Estimation



$Y'$ $X'$

camera coordinate

$Z'$

$C$

Camera

Image

3D rotation $\mathbf{R}$

$Z$

$Y$

$O$

$X$

object coordinate

3D Translation
$\mathbf{T} = (T_x, T_y, T_z)^T$

3D models

3D world

11

# Traditional Methods for 6D Object Pose Estimation

- Feature matching-based methods
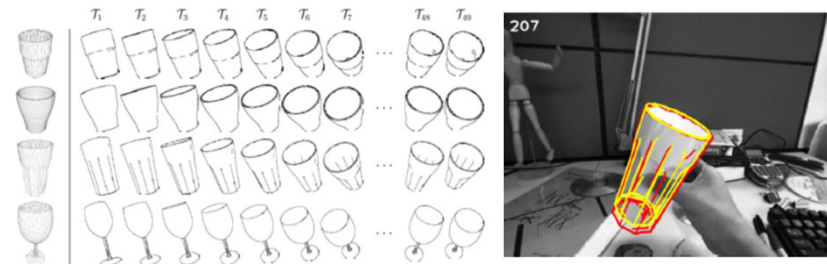


Rothganger-Lazebnik-Schmid-Ponce, IJCV'06



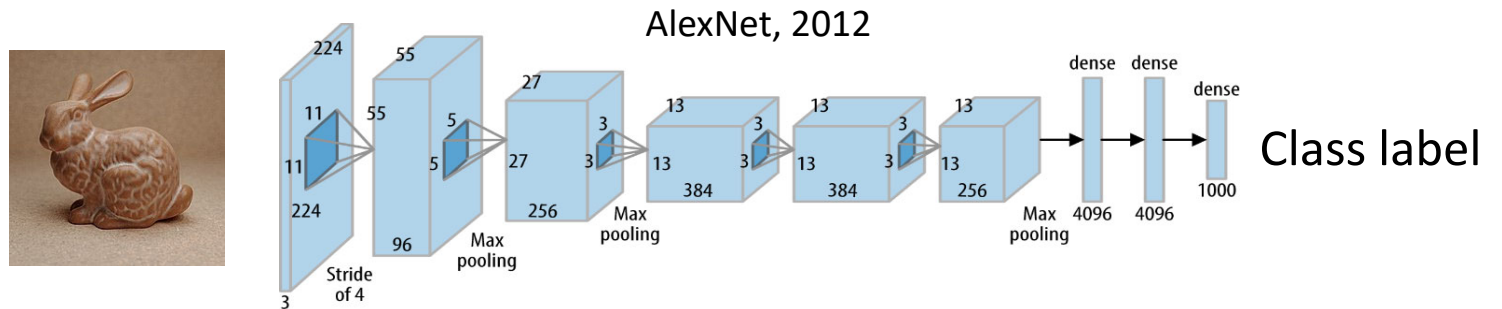Collet-Martinez-Srinivasa, IJRR'11

- Template matching-based methods



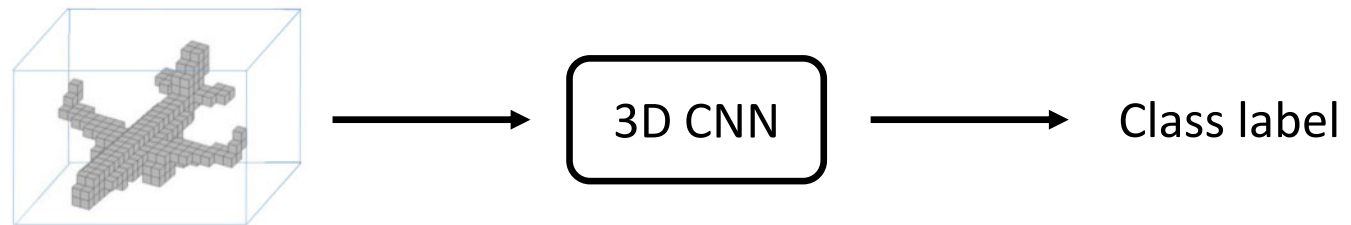Hinterstoisser-Lepetit-Ilic-Holzer-Bradski-Konolige-Navab, ACCV'12



Choi-Christensen, IROS'12
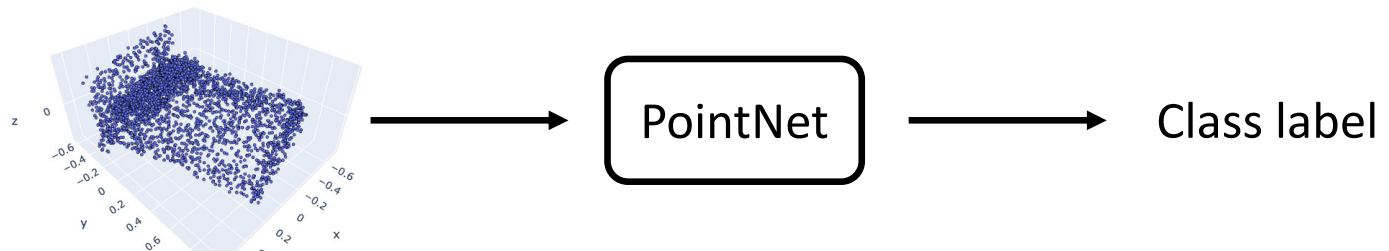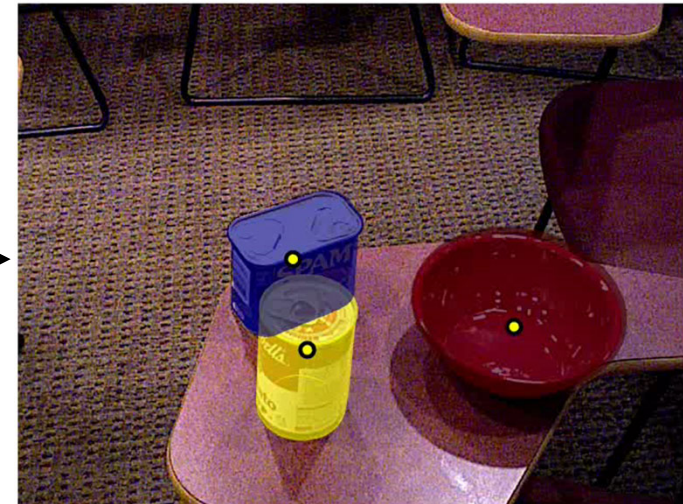
12

# Deep Learning for Visual Recognition

- Images

AlexNet, 2012

Class label

- Voxels

3D CNN → Class label

- Point Clouds

PointNet → Class label

# PoseCNN: the First End-to-end 6D Pose Estimation Network



PoseCNN

✓ Texture-less objects
✓ Symmetric objects
✓ Occlusions

**Xiang**-Schmidt-Narayanan-Fox, RSS'18

# PoseCNN: the First End-to-end 6D Pose Estimation Network



Semantic segmentation

Center direction X

Center direction Y

Center distance

$\mathbf{T} = (T_x, T_y, T_z)^T$

Input image

Feature extraction

3D translation estimation

Hough voting layer

For each RoI

$\mathbf{R}$

3D Rotation regression

6D Poses

**Xiang**-Schmidt-Narayanan-Fox, RSS'18

15

# PoseCNN: the First End-to-end 6D Pose Estimation Network



Segmentation and Detection

Poses

3D World

Input image

**Xiang**-Schmidt-Narayanan-Fox, RSS'18

# The Sim-to-Real Gap

## Synthetic images



Training → PoseCNN

Domain randomization

## Lighting and background



## Texture



## Moving Part



17

# Self-supervised 6D Object Pose Estimation
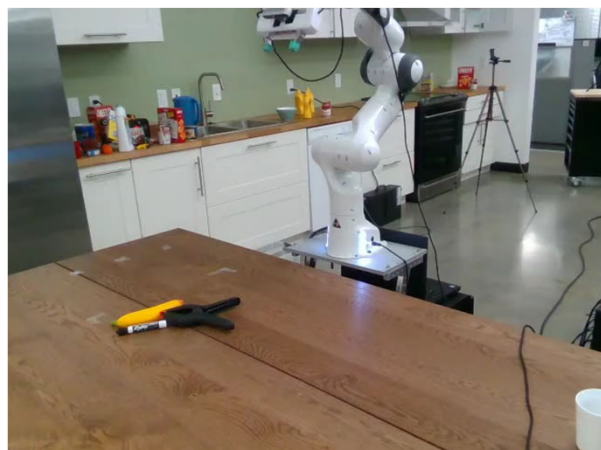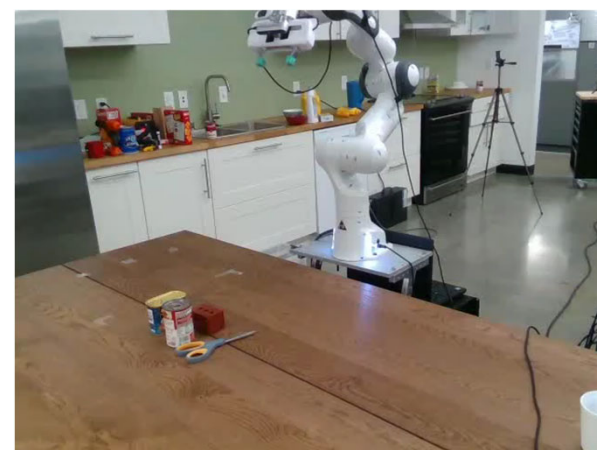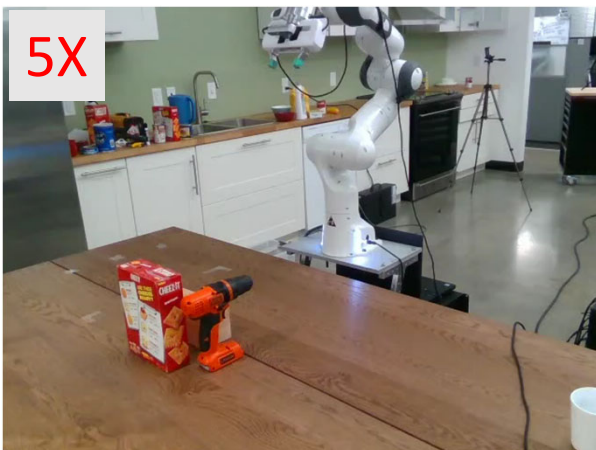
Interactive real-world data collection

Generated pose annotations
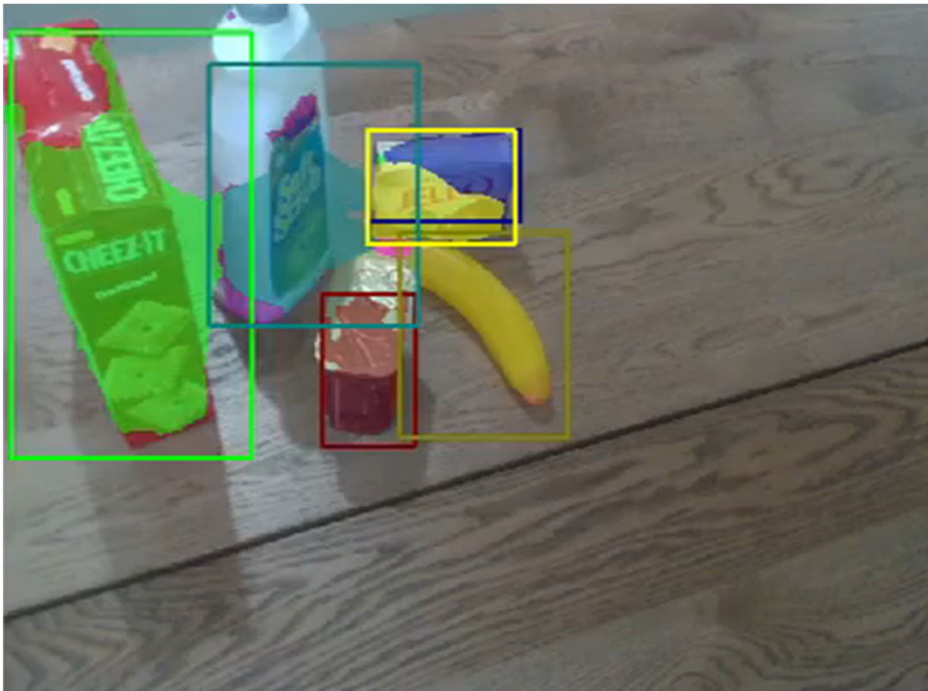


Overlay of rendering onto image

Deng-**Xiang**-Mousavian-Eppner-Bretl-Fox, ICRA'20

# Self-supervised 6D Object Pose Estimation

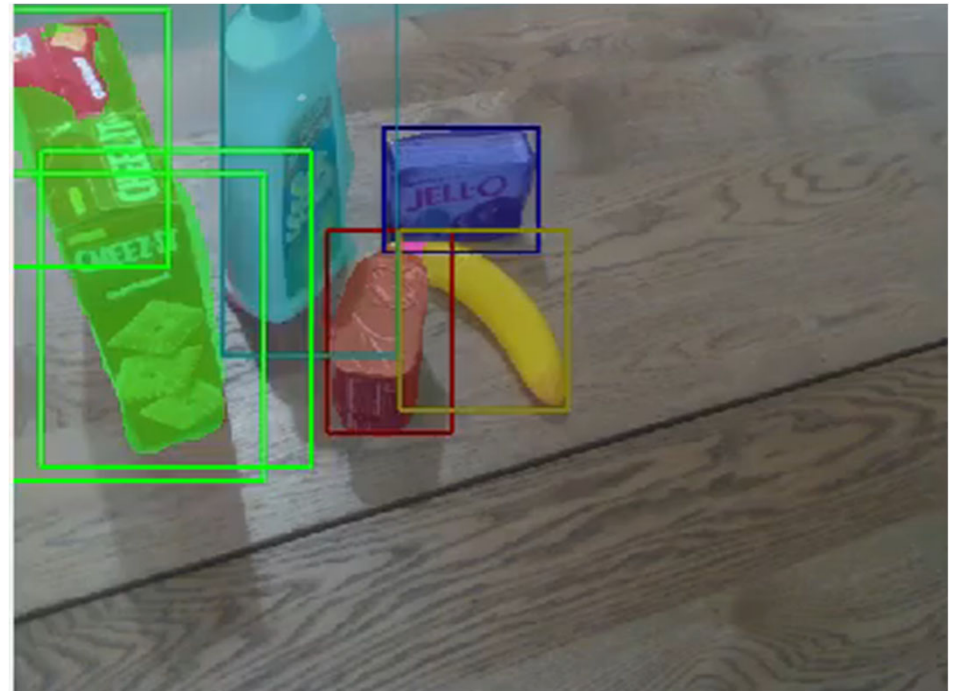12 robot hours, 497 scenes
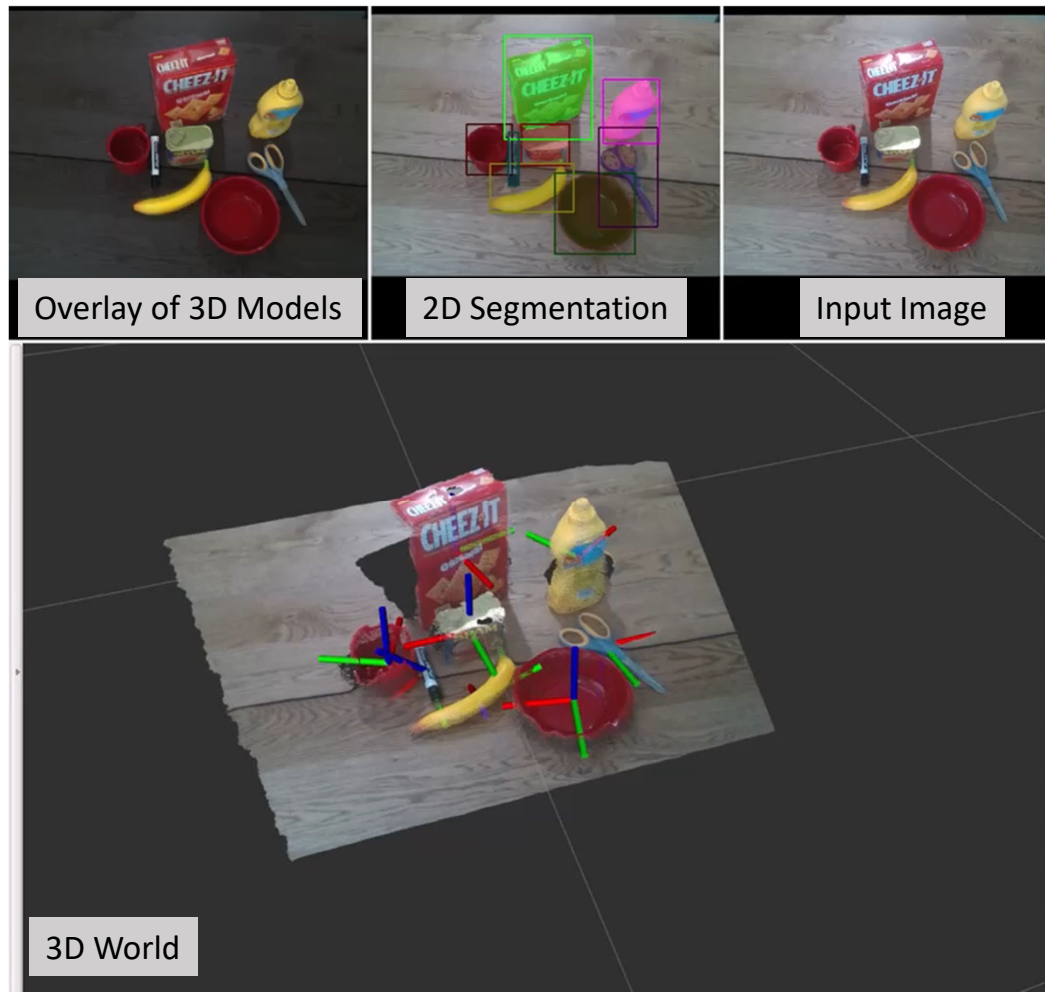6,541 RGB-D images,
22,851 object instances

5X

Deng-**Xiang**-Mousavian-Eppner-Bretl-Fox, ICRA'20

19

# Self-supervised 6D Object Pose Estimation

PoseCNN
trained with only synthetic data

PoseCNN
fine-tuned with self-annotated data



Deng-**Xiang**-Mousavian-Eppner-Bretl-Fox, ICRA'20

20

# Perception: Model-based 6D Object Pose Estimation



Overlay of 3D Models | 2D Segmentation | Input Image

3D World



3D models

PoseCNN: **Xiang**-Schmidt-Narayanan-Fox, RSS'18
DeepIM: Li-Wang-Ji-**Xiang**-Fox, ECCV'18 Oral, IJCV'19
PoseRBPF: Deng-Mousavian-**Xiang**-Xia-Bretl-Fox, RSS'19, T-RO'21
Self-supervision 6D Pose: Deng-**Xiang**-Mousavian-Eppner-Bretl-Fox, ICRA'20

Codes available online
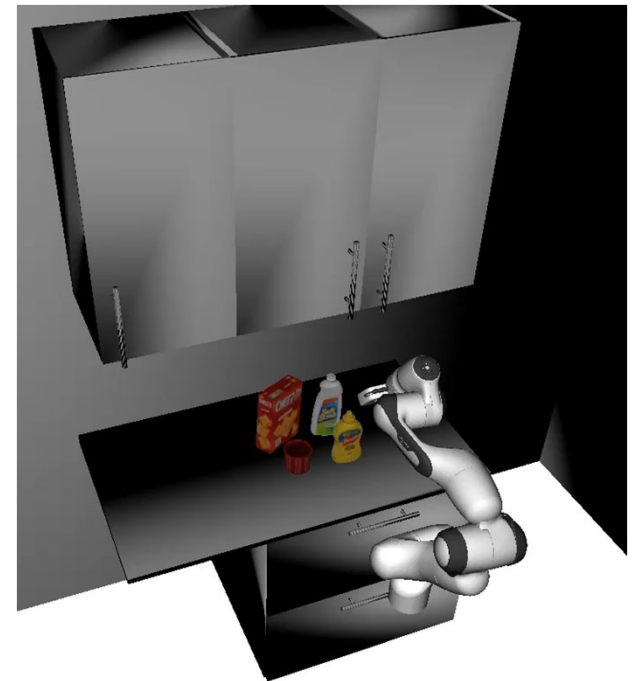
21

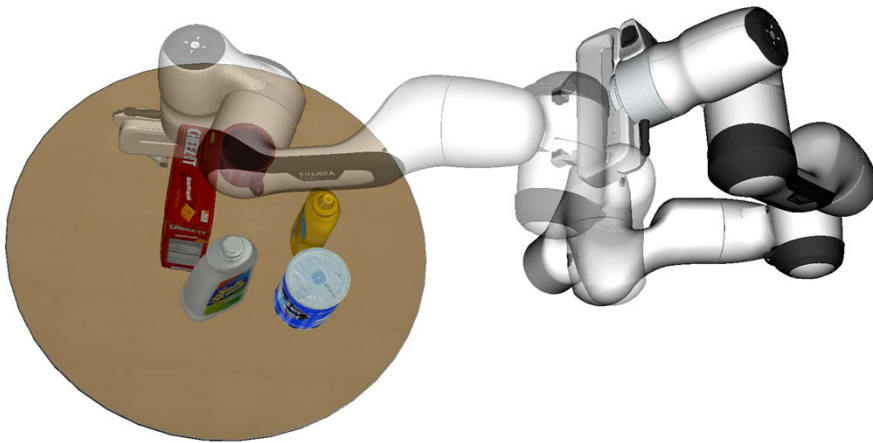# Manipulation Planning

Input image

6D Object Pose Estimation

3D models

Planning scene



22

# Manipulation Planning

## Arm Motion Planning



We need to specify a goal configuration.

## Grasp Planning



No arm motion is considered.

**Sampling-based methods**:
PRM: Kavraki-Svestka-Latombe-Overmars, T-RA'96
RRT: LaValle, Techincal Report'98
RRT-Connect: Kuffner-LaValle, ICRA'00
SMRM: Alterovitz-Simeon-Goldberg, RSS'07
RRT*: Karaman-Frazzoli, IJRR'11
FMT: Janson-Schmerling-Clark-Pavone, IJRR'15

**Trajectory optimization**:
CHOMP: Ratliff-Zucker-Bagnell-Srinivasa, ICRA'09
STOMP: Kalakrishnan-Chitta-Theodorou-Pastor-Schaal, ICRA'11
TrajOpt: Schulman-Duan-Ho-Lee-Awwal-Bradlow-Pan-Patil-Goldberg-Abbeel, IJRR'14
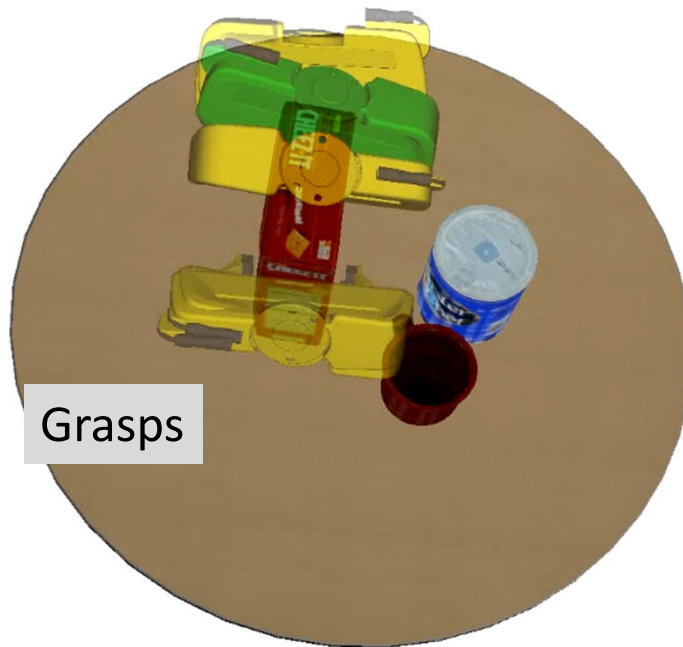GPMP2: Mukadam-Dong-Yan-Dellaert-Boots, IJRR'18

Nguyen, IJRR'88
Ferrari-Canny, ICRA'92
Chen-Burdick, T-RA'93
Graspit!: Miller-Allen, RA Magazine'04
Ciocarlie-Goldfeder-Allen, RSS Workshop'07
ten Pas-Gualtieri-Saenko-Platt, IJRR'17
Fan-Lin-Tang-Tomizuka, CASE'18
Mousavian-Eppner-Fox, ICCV'19

23

# OMG Planner: An Optimization-based Motion and Grasp Planner



Grasps

Arm motion

Joint Motion and Grasp Planning
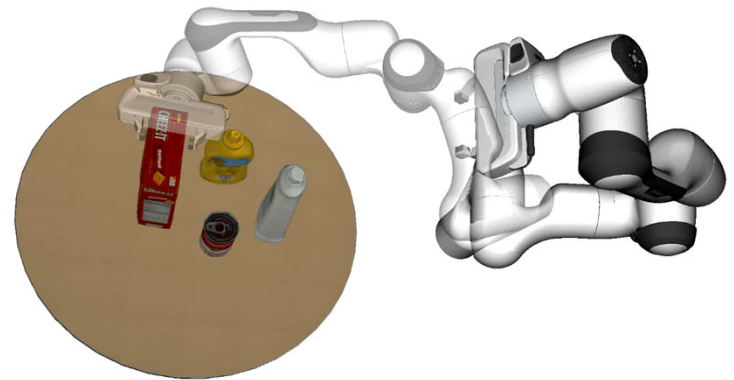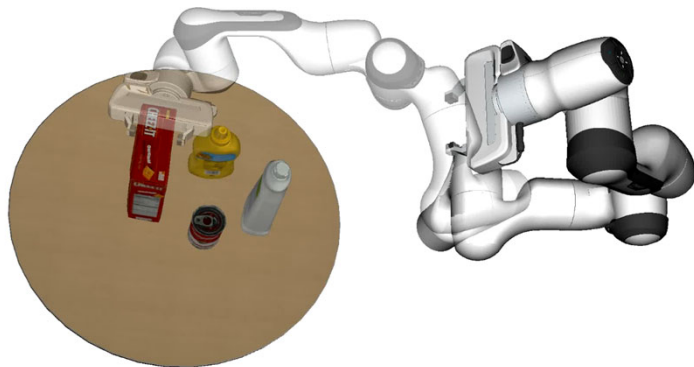
# Trajectory Optimization: CHOMP

$$f_{\mathrm{motion}}(\xi) = f_{\mathrm{obstacle}}(\xi) + \lambda f_{\mathrm{smooth}}(\xi)$$
$$\xi = (q_1, \ldots, q_T)$$ A trajectory of robot joint configurations
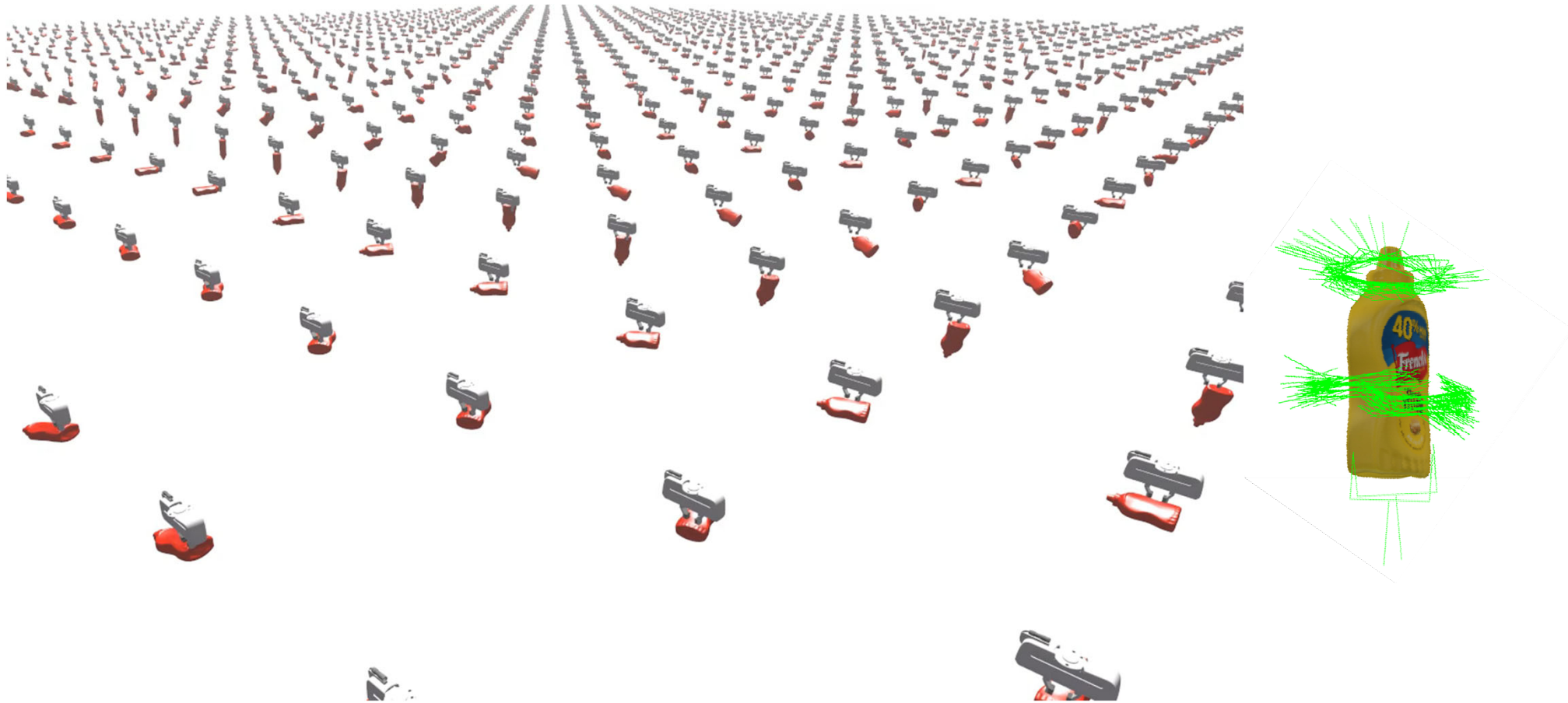
N steps gradient descent

Initial trajectory with collision

Final trajectory



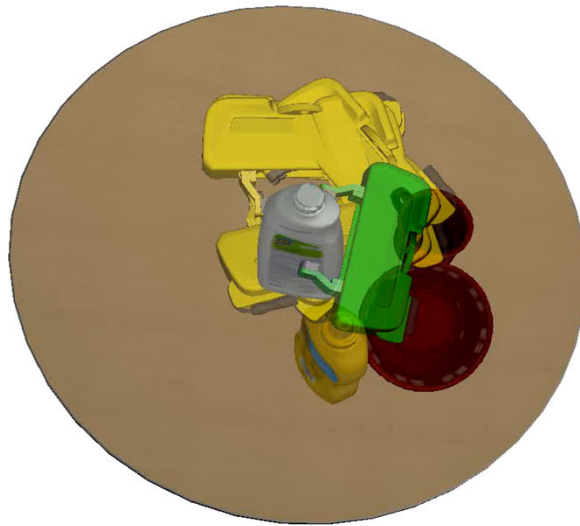Covariant Hamiltonian Optimization for Motion Planning (CHOMP): Ratliff-Zucker-Bagnell-Srinivasa, ICRA'09
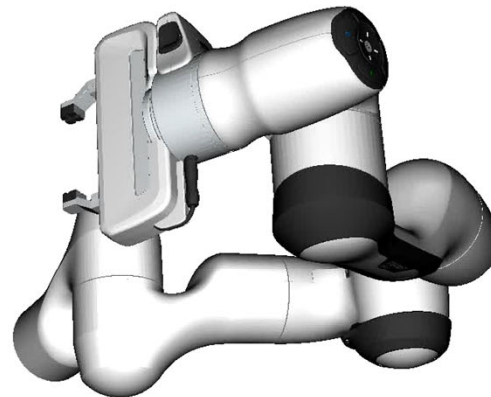
# Grasp Planning: A Physics-based Approach

# OMG Planner: Trajectory Optimization and Grasp Selection
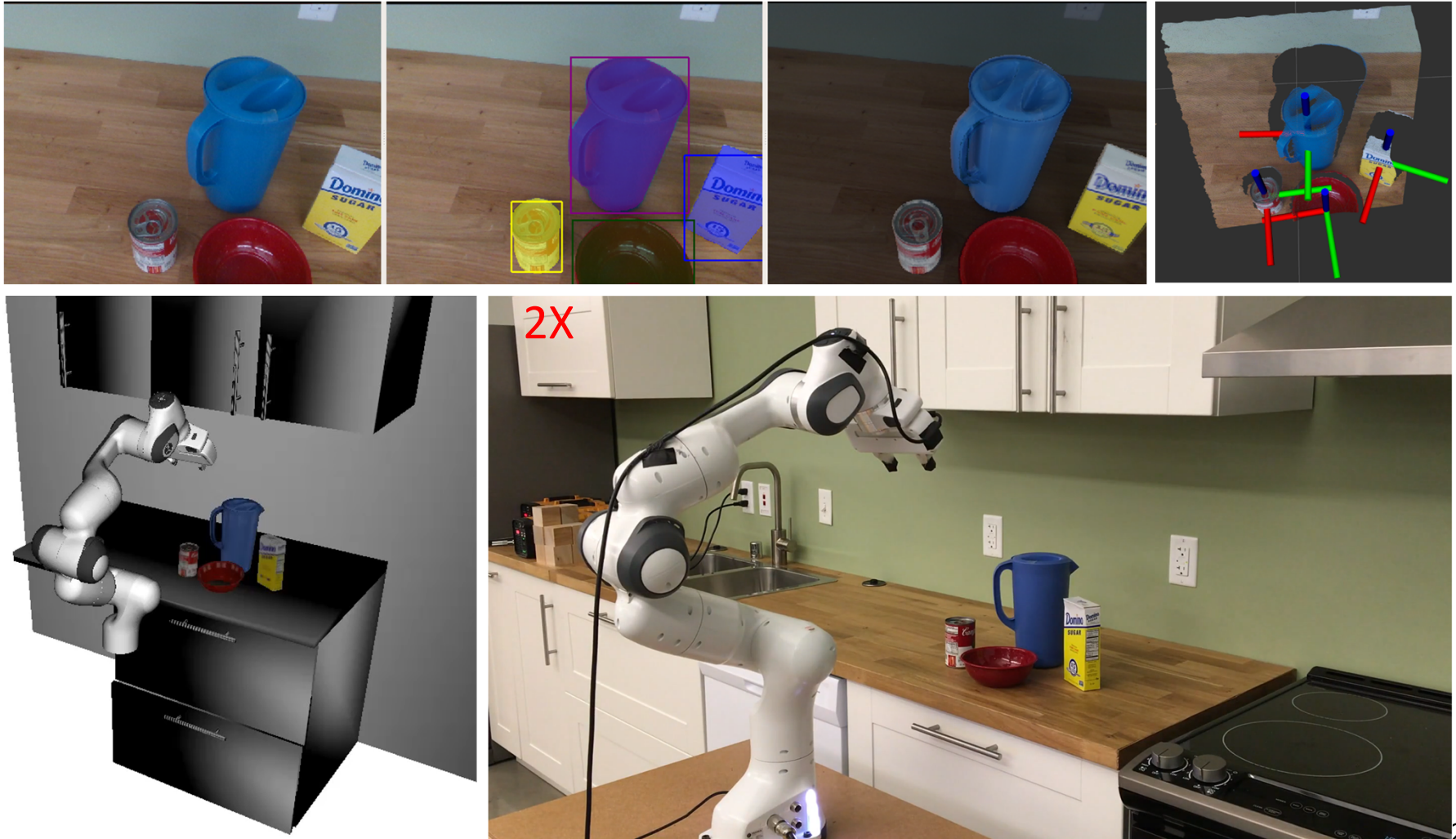
OMG Iter: 50



100 grasps

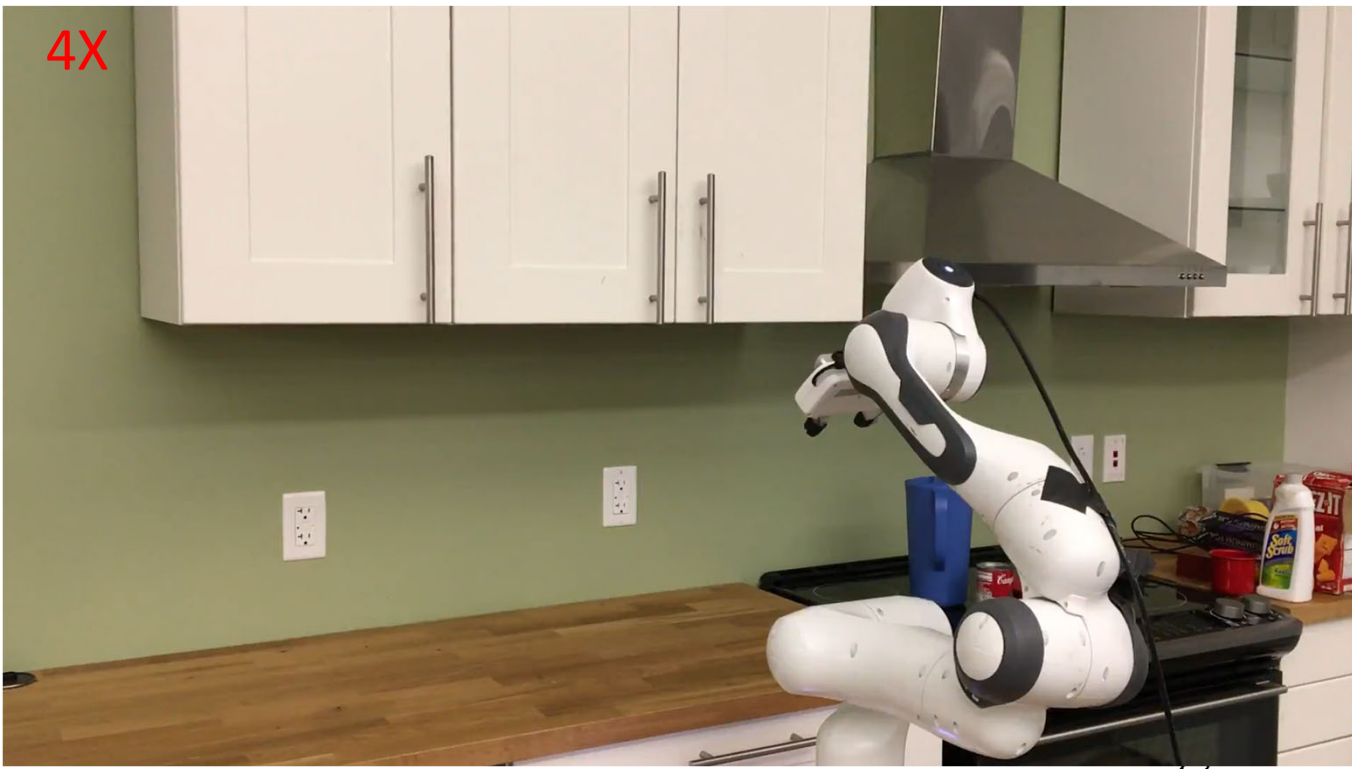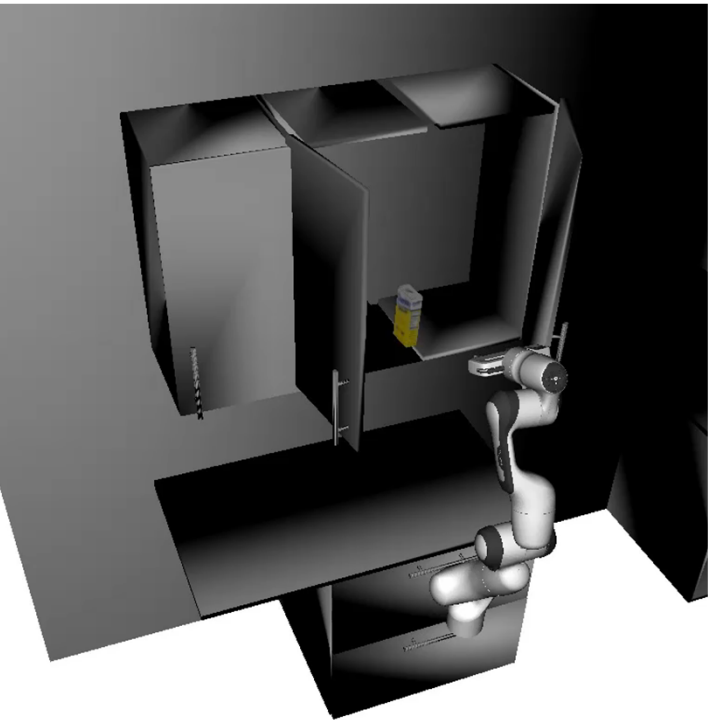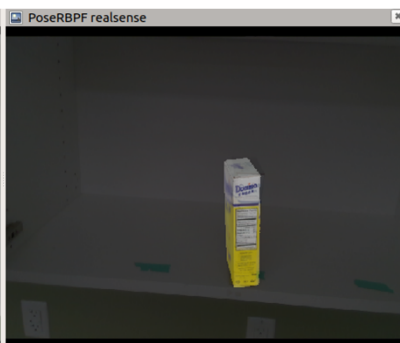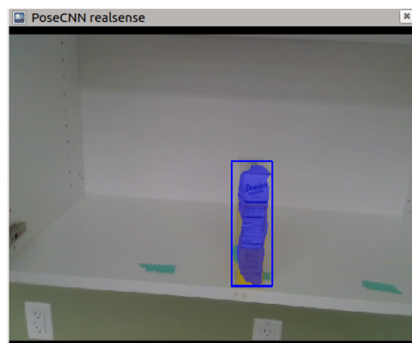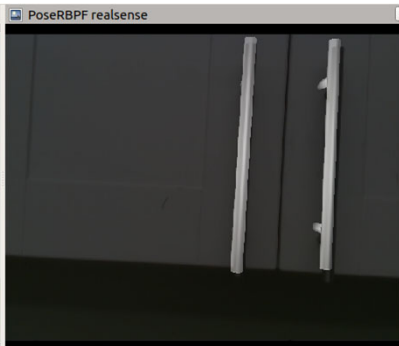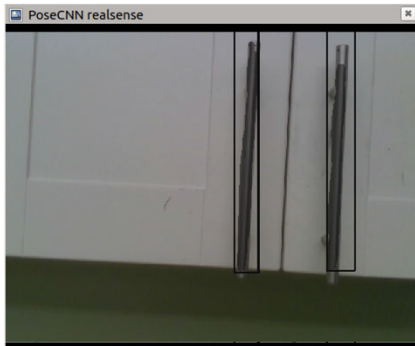Modeling the goal set distribution
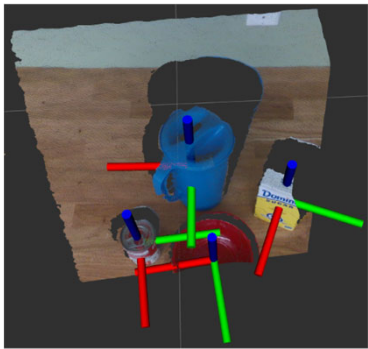
Code available online

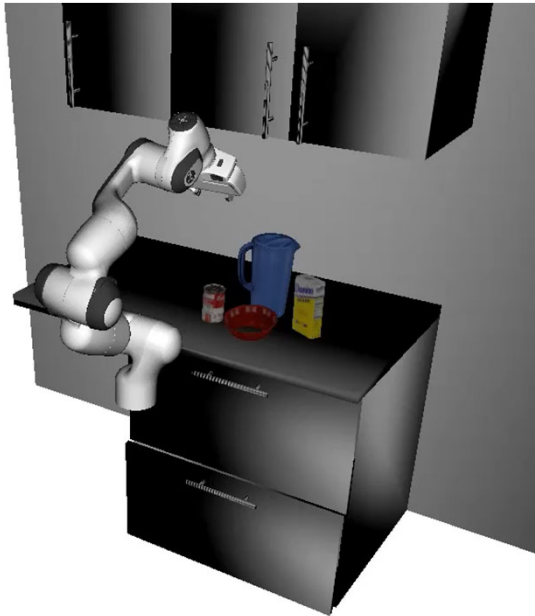# Real-world Manipulation with 6D Pose Estimation and Planning

# Model-based Robot Manipulation

6D Object Pose Estimation
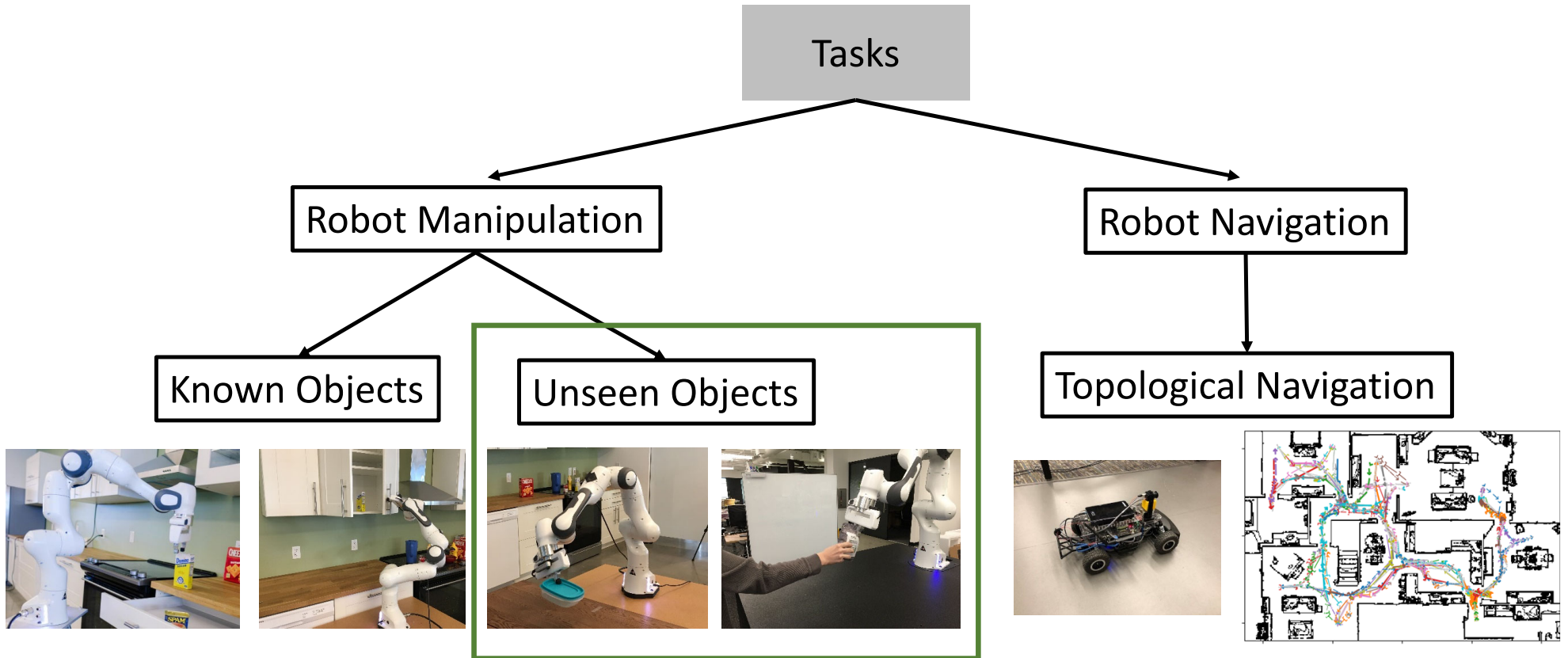
Motion and Grasp Planning



We need to have 3D models of objects

How can we enable robots to manipulate unseen objects?

# Outline



Tasks

Robot Manipulation

Robot Navigation

Known Objects

Unseen Objects

Topological Navigation

# Model-free Robot Manipulation
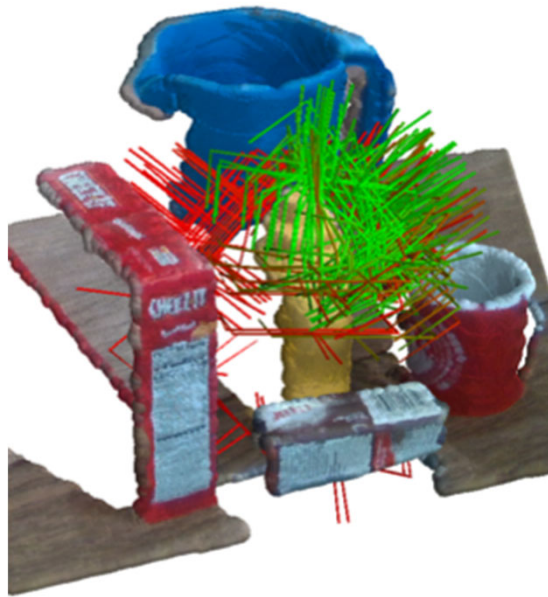


Perception → Planning → Control

Unseen object instance segmentation

Grasp planning from point clouds

Position control to reach grasp

Figure Credit: Murali-Mousavian-Eppner-Paxton-Fox, ICRA'20

32

# Perception: Unseen Object Instance Segmentation
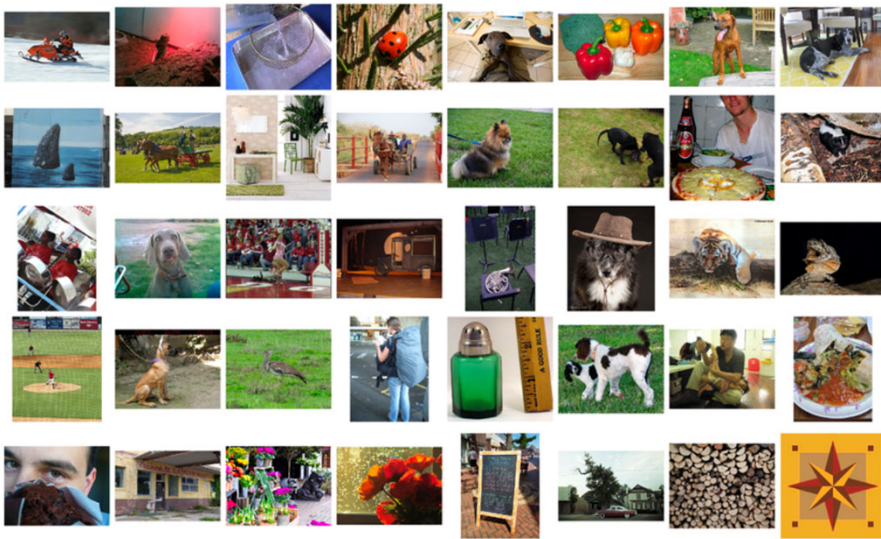


Xie-**Xiang**-Mousavian-Fox, CoRL'19, T-RO'21

Codes available online

**Xiang**-Xie-Mousavian-Fox, CoRL'20

Training on synthetic data, transferring well to the real images for segmenting unseen objects

# Learning the Concept of "Objects"
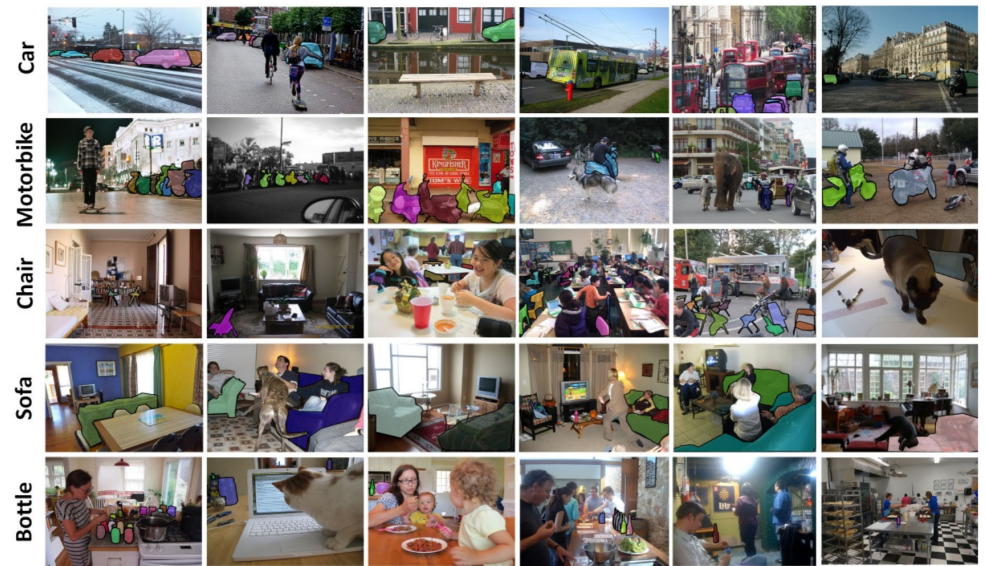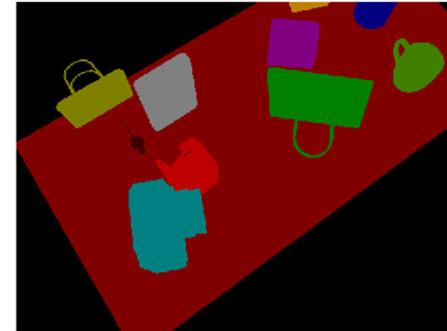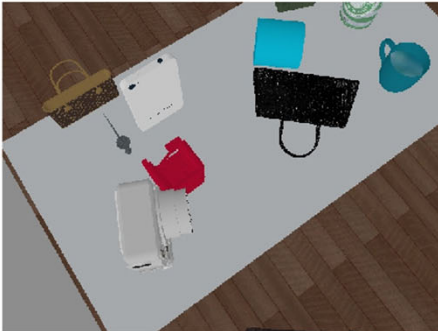
- Learning from data



ImageNet: Deng-Dong-Socher-Li-Li-Fei-Fei, CVPR'09



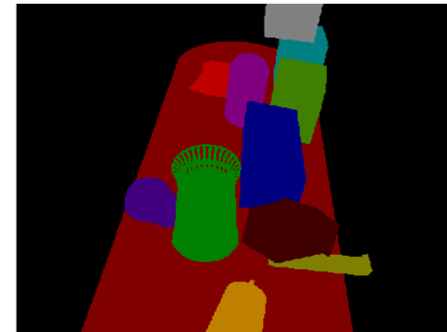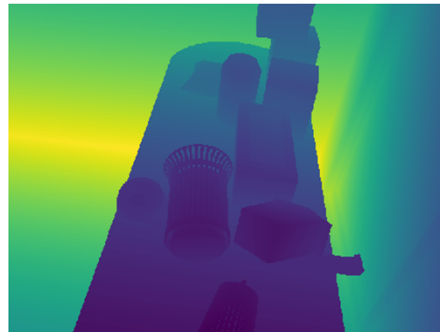COCO: Lin-Maire-Belongie-Bourdev-Girshick-Hays-Perona-Ramanan-Zitnick-Dollar, ECCV'14

Internet Images, not suitable for indoor robotic settings

34

# Learning from Synthetic Data



RGB       Depth       Instance Label

ShapeNet objects in the PyBullet simulator

40,000 scenes
7 RGB-D images per scene

Need to deal with the sim-to-real gap

Tabletop Object Dataset: Xie-**Xiang**-Mousavian-Fox, CoRL'19

35

# Unseen Object Instance Segmentation: Learning RGB-D Feature Embeddings



Instance Label for Training

Dense Feature Map

Metric Learning Loss

Fully Convolutional Network

RGB

Depth

- ● Sampled feature
- ✖ Cluster center
- → Intra-cluster
- ↔ Inter-cluster

**Xiang**-Xie-Mousavian-Fox, CoRL'20

36

Input Image

Feature Map

Output Label

**Xiang**-Xie-Mousavian-Fox, CoRL'20

37

# Grasp Planning from Partially Observed Point Clouds



Input Image

Object Point Cloud

**Grasp Sampler**

Sampled Grasps

**Grasp Evaluator**

Assessed Grasps

**Grasp Refinement**

6-DOF GraspNet: Mousavian-Eppner-Fox, ICCV'19

38

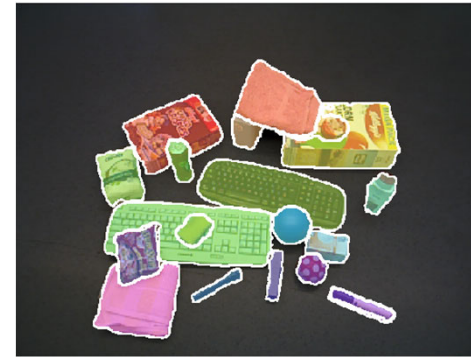# Grasping Unseen Objects



Unseen Object Instance Segmentation:
Xie-**Xiang**-Mousavian-Fox, CoRL'19, T-RO'21
**Xiang**-Xie-Mousavian-Fox, CoRL'20

6-DOF GraspNet:
Mousavian-Eppner-Fox, ICCV'19

# Open-Loop VS. Closed-Loop



Tasks

Perception → Planning → Control

Sensing

World

Action

# Closed-loop Robot Control with Markov Decision Processes



State $s_t$

Reward $r_t$

Robot

Action $a_t$

Environment

$s_{t+1}$

Reinforcement Learning:
Imitation Learning: $a_t = \pi(s_t)$

# Learning Closed-Loop Control Polices for 6D Grasping



Segmentation

Image

Point cloud

State $s_t$

Policy

Deep Neural Network

Action $a_t$

Closed-Loop

3D Translation
3D Rotation

Perception

Control

Wang-**Xiang**-Fox, in arXiv'21

No planning?

42

# Learning from Demonstration with the OMG-Planner

50,000 trajectories
1,500 3D shapes



Wang-**Xiang**-Fox, in arXiv'21

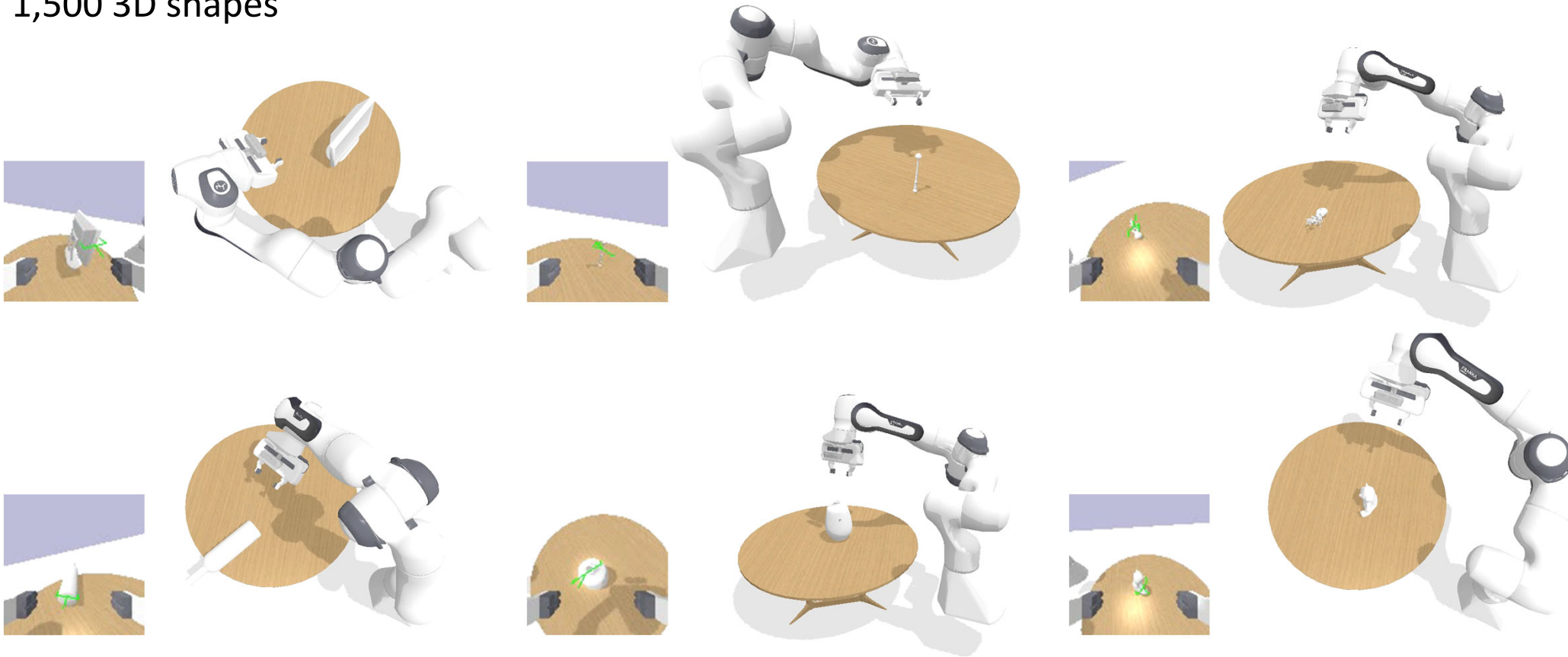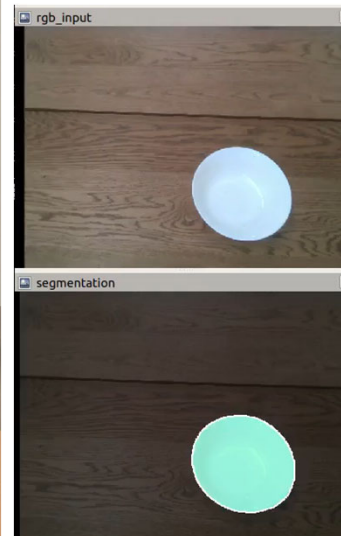# Our Learned Policy in the Real World

# Closed-Loop Human-Robot Handover



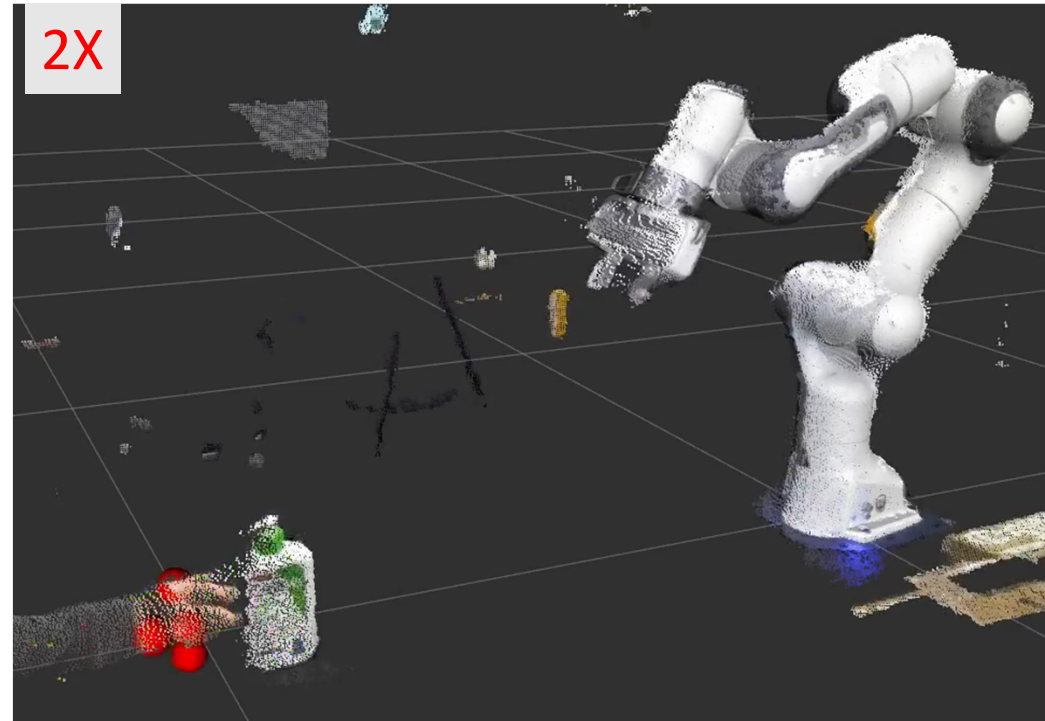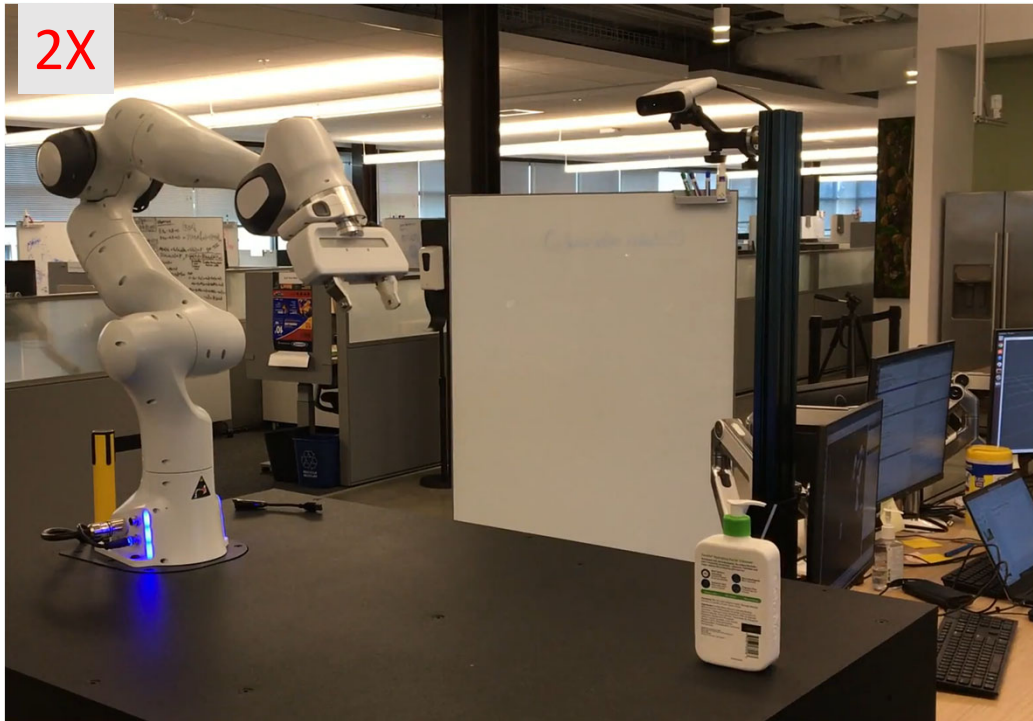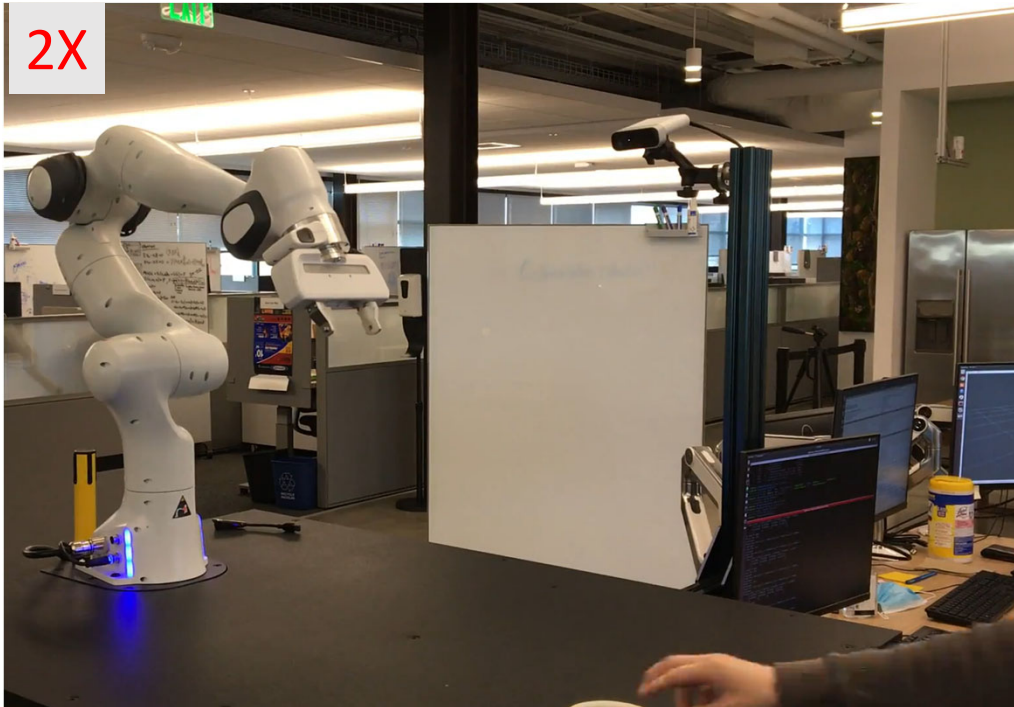Yang-Paxton-Mousavian-Chao-Cakmak-Fox, in arXiv'20
Wang-**Xiang**-Fox, in arXiv'21

45

# Closed-Loop Human-Robot Handover



Yang-Paxton-Mousavian-Chao-Cakmak-Fox, in arXiv'20
Wang-**Xiang**-Fox, in arXiv'21

# Manipulation and Navigation

# Outline



Tasks

Robot Manipulation

Robot Navigation

Known Objects

Unseen Objects

Topological Navigation

# Traditional Robot Navigation

| Perception | → | Planning | → | Control |
|---|---|---|---|---|

Simultaneous localization and mapping (SLAM)

Path planning

Path following



Laser-based SLAM
2D occupancy grid map

Limitations of SLAM-based navigation
- 3D reconstruction is expensive
- Detailed 3D geometry information may not be necessary

49

# Topological Navigation

Dense Trajectories

Sparse Topological Map

Reachability
Estimator

Local Controller
A learned neural network
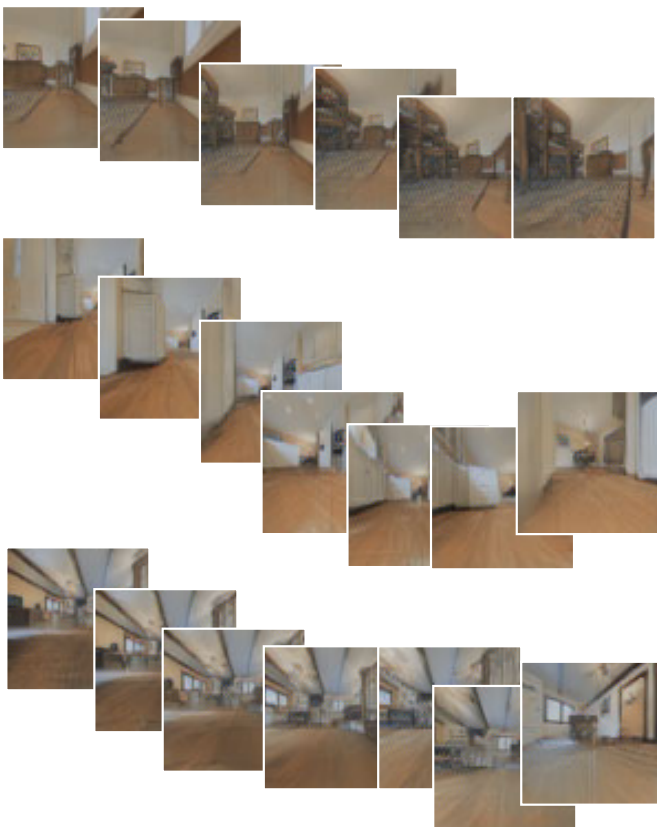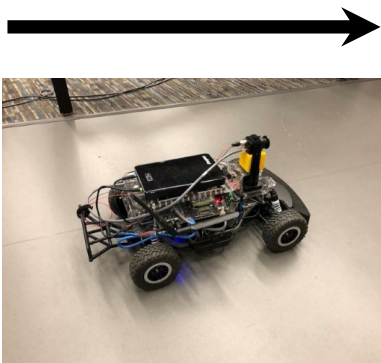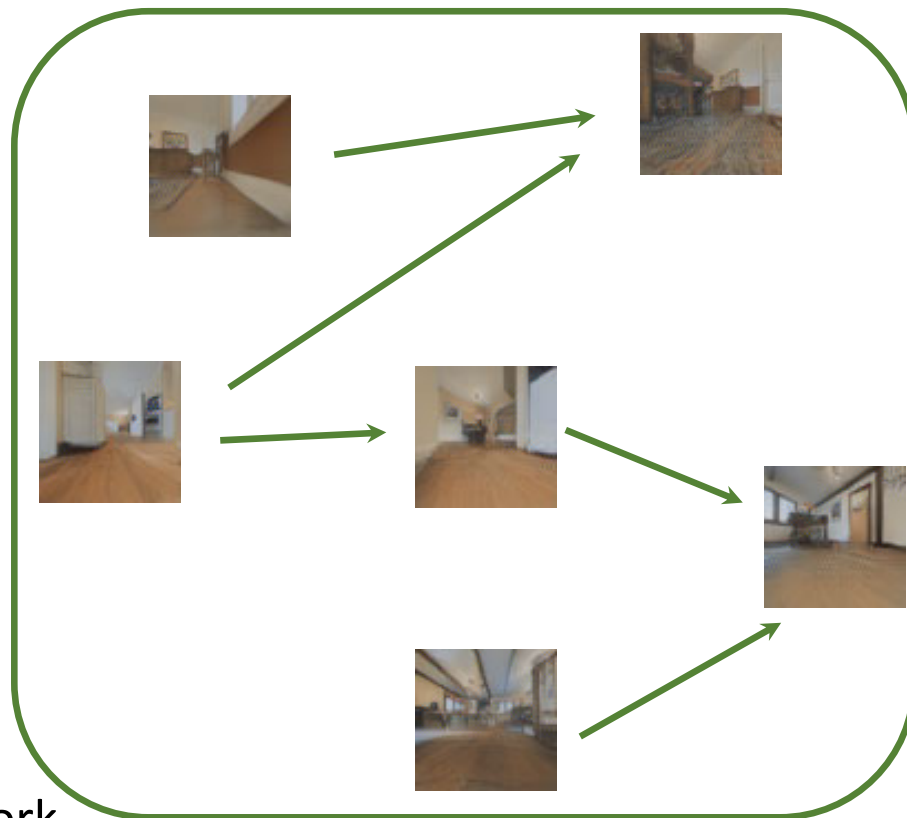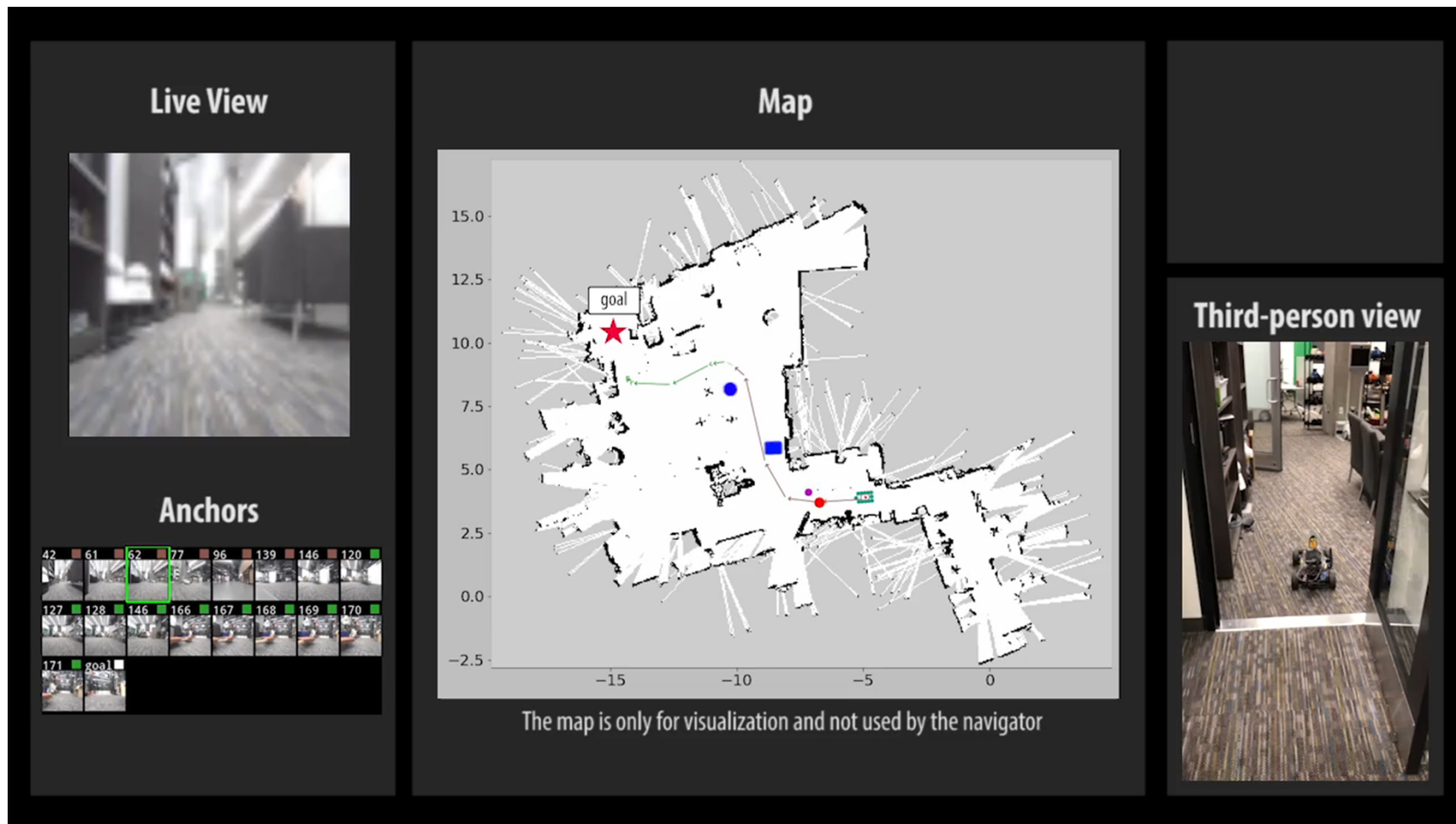
# Topological Navigation

Meng-Ratliff-**Xiang**-Fox, ICRA'19, '20
Meng-**Xiang**-Fox, RA-L'21



51

# Summary



```
                        ┌─────────┐
                        │  Tasks  │
                        └─────────┘
              ┌──────────────┴──────────────┐
    ┌───────────────────┐         ┌────────────────────┐
    │ Robot Manipulation │         │  Robot Navigation  │
    └───────────────────┘         └────────────────────┘
       ┌──────┴───────┐                    │
┌──────────────┐ ┌──────────────┐  ┌──────────────────────┐
│ Known Objects │ │ Unseen Objects │  │ Topological Navigation │
└──────────────┘ └──────────────┘  └──────────────────────┘
```

| Perception | Planning | Control | Learning |

# Future Work: Long-horizon Tasks in Human Environments



Multiple Tasks

Task Diversity

Generalizability

**Intelligent Robots**
- Make a cup of coffee
- Set a dining table
- Execute human instructions
  - "Bring a bottle of water"

- Manipulation
- Navigation

**Industrial robots**

Single Task

Single Environment     Environment Diversity     Multiple Environments

# Future Work: Learning Robot Skills and Building Robotic Systems

## Robot Skills Generalizable and Shareable

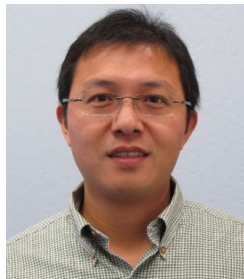| Perception | Planning | Control | Learning |
|---|---|---|---|
| • Understand objects, scenes and space<br>• Understand humans and language | • Task planning<br>• Motion planning | • Learning task-specific controllers | • Supervised Learning<br>• Imitation Learning<br>• Reinforcement Learning |

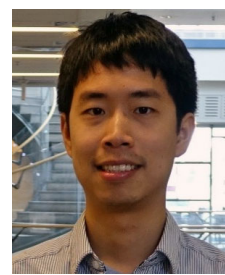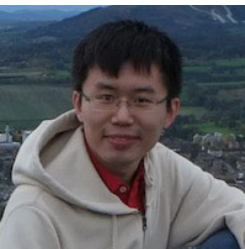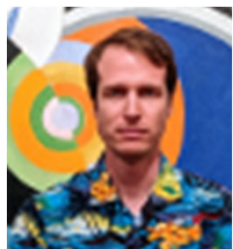**Deploy** ⬇  ⬆ **Improve**

## Robotic Systems



- Closing the perception, planning and control loop
- Self-supervised learning
- Life-long learning

54

# Our Missions of the Future Research Lab

- Advancing robot perception, planning and control

- Building intelligent robotic systems

- Open-sourcing and sharing

- Collaborating

# Acknowledgements

## Thank you!