# Learning How-To Knowledge from the Web

Yuke Zhu

IROS 2019

Stanford University

NVIDIA

TEXAS
The University of Texas at Austin

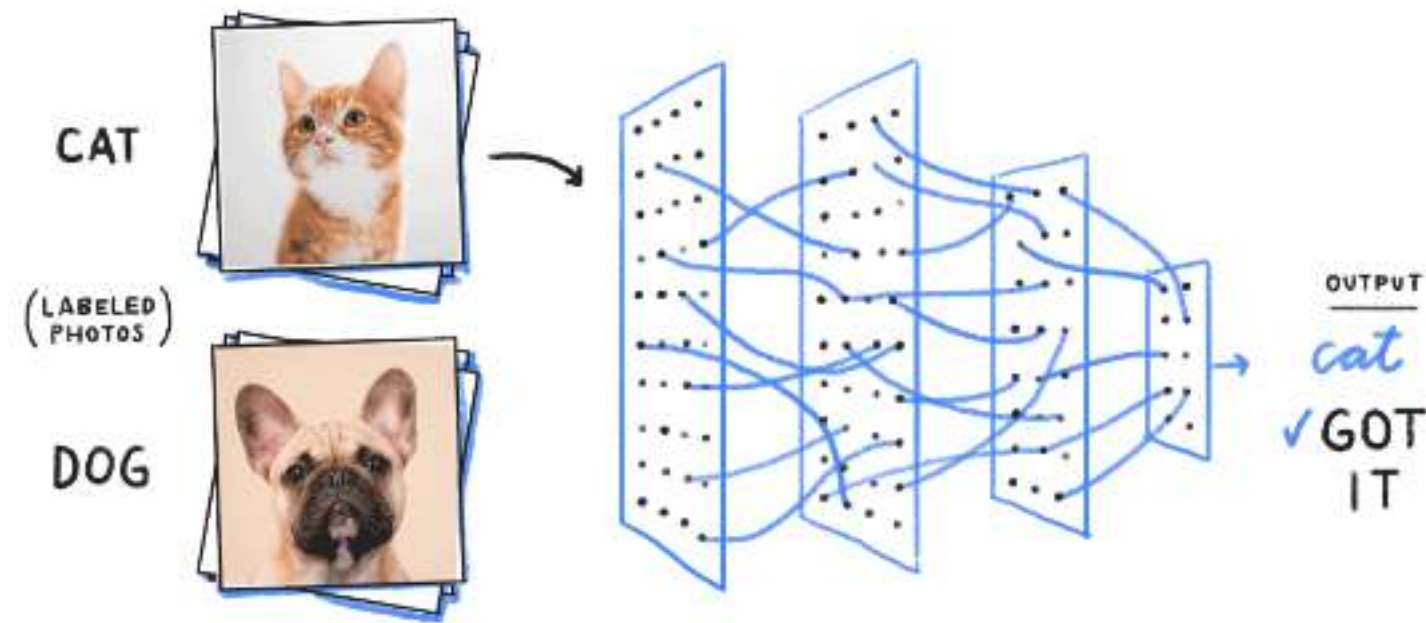# Advances in Artificial Intelligence



Visual Recognition

Machine Translation

Question Answering

# The Unsung Hero: Web Data

## Visual Recognition

**ImageNet**
[Deng et al. 2009]

14 million web images
annotated by AMT workers

## Machine Translation

**Google NMT**
[Wu et al. 2016]

WMT En→Fr dataset
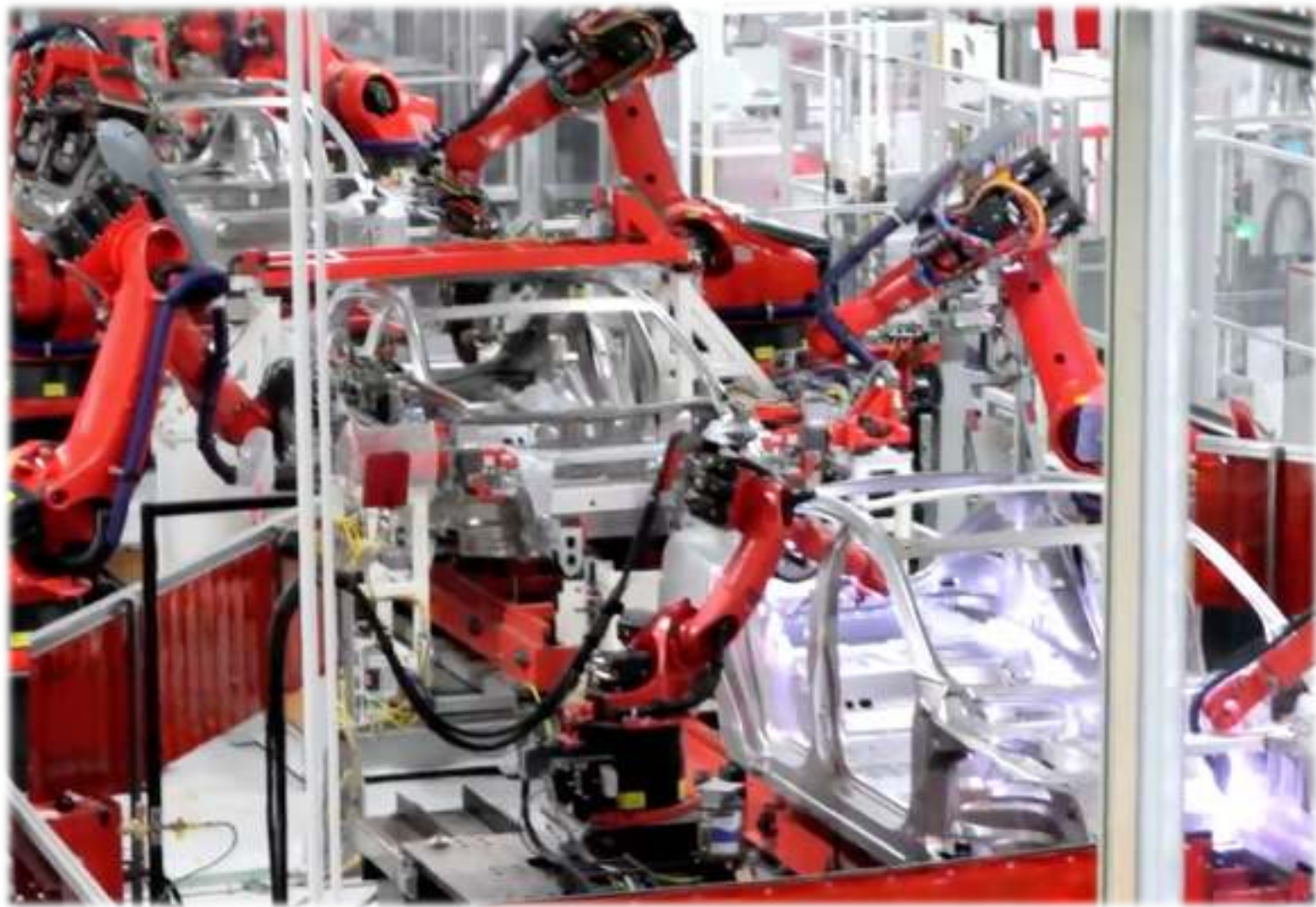with 36M sentence pairs

## Question Answering

**SQuAD QA Dataset**
[Rajpurkar et al. 2016]

100,000+ questions posed
by crowdworkers on a set
of Wikipedia articles

# The Unsung Hero: Web Data

What's the role of web data in improving robot intelligence?



?

Traditional form of automation                    Intelligent robots in real world
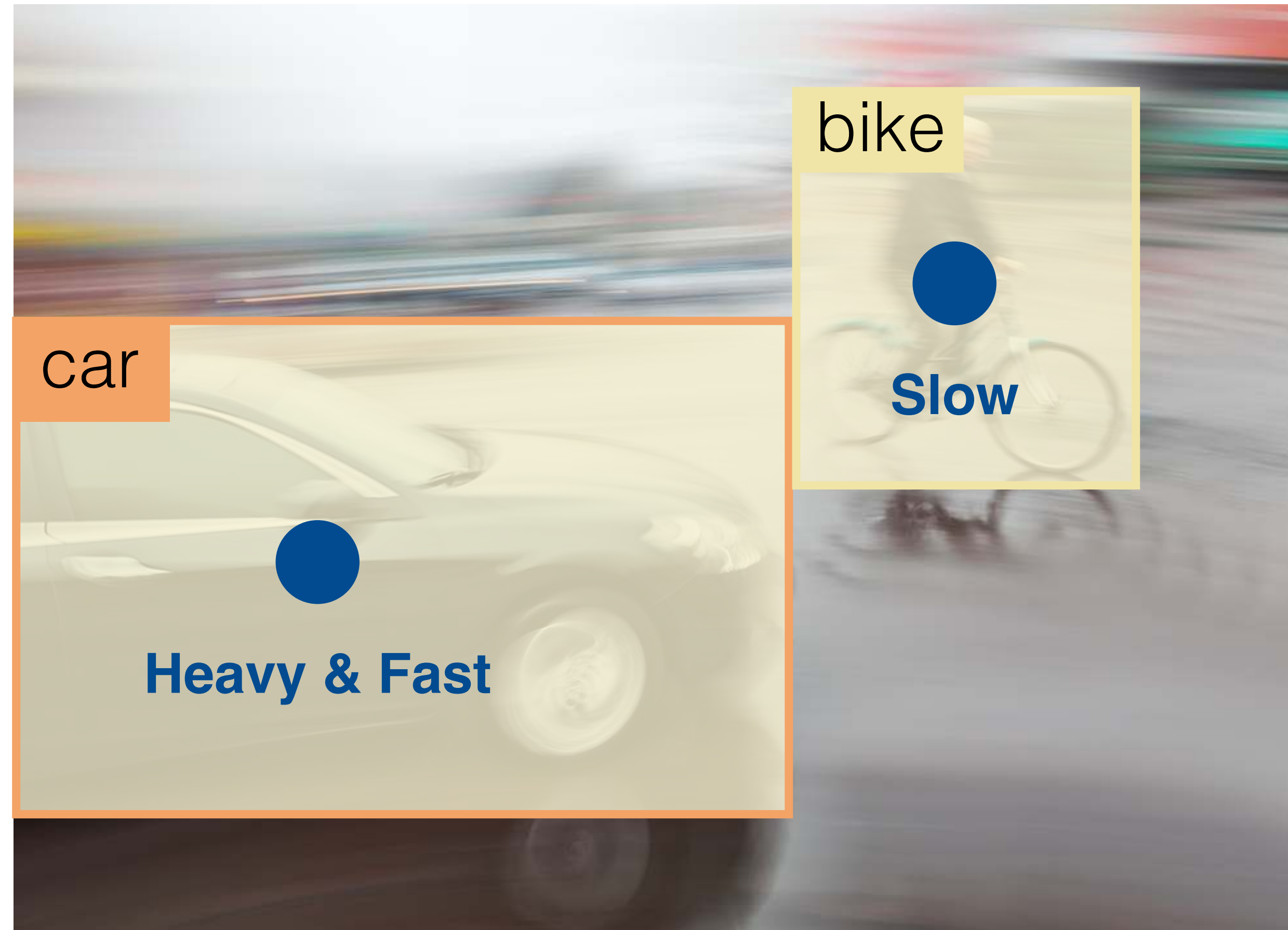
# What knowledge do we need for robotics?



"To accelerate or to brake?"

# What knowledge do we need for robotics?

**Declarative knowledge**

Understanding the world

❖ Describes facts
   of the world

❖ Easy to articulate
   (conscious)



bike

car

**Slow**

**Heavy & Fast**

Knowledge of "That-Is"

# What knowledge do we need for robotics?

**Declarative knowledge**

Understanding the world

❖ Describes facts
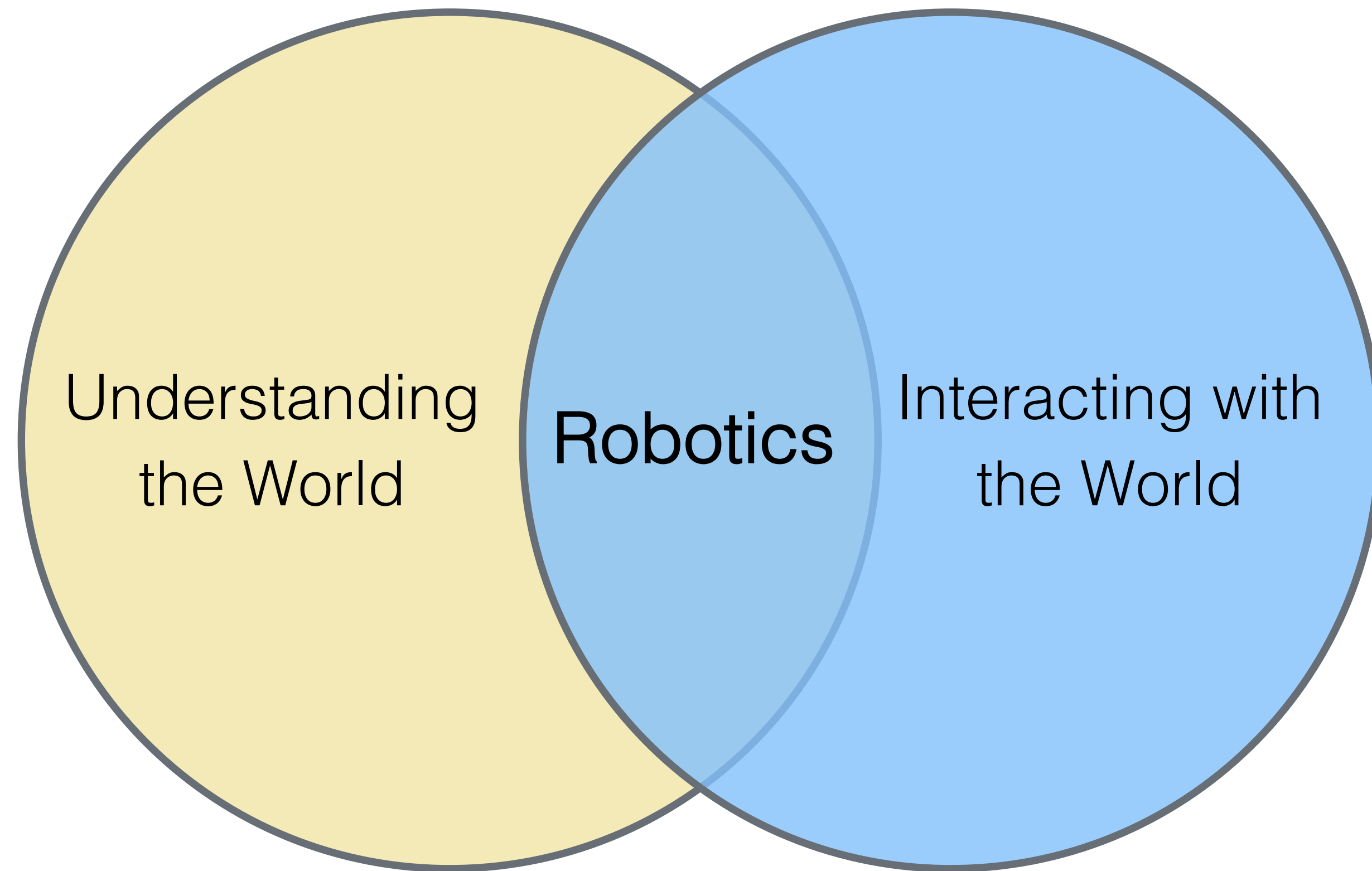   of the world

❖ Easy to articulate
   (conscious)

**Procedural knowledge**

Interacting with the world

❖ Describes **how to**
   perform tasks
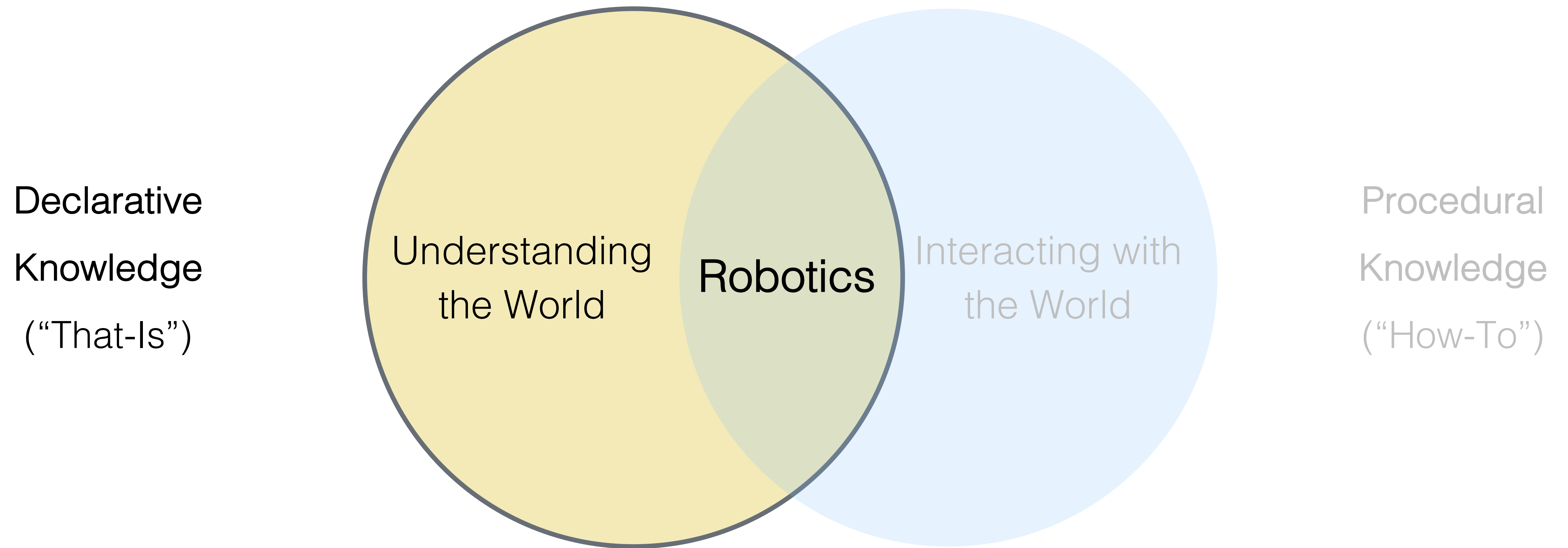
❖ Hard to pinpoint
   (unconscious)

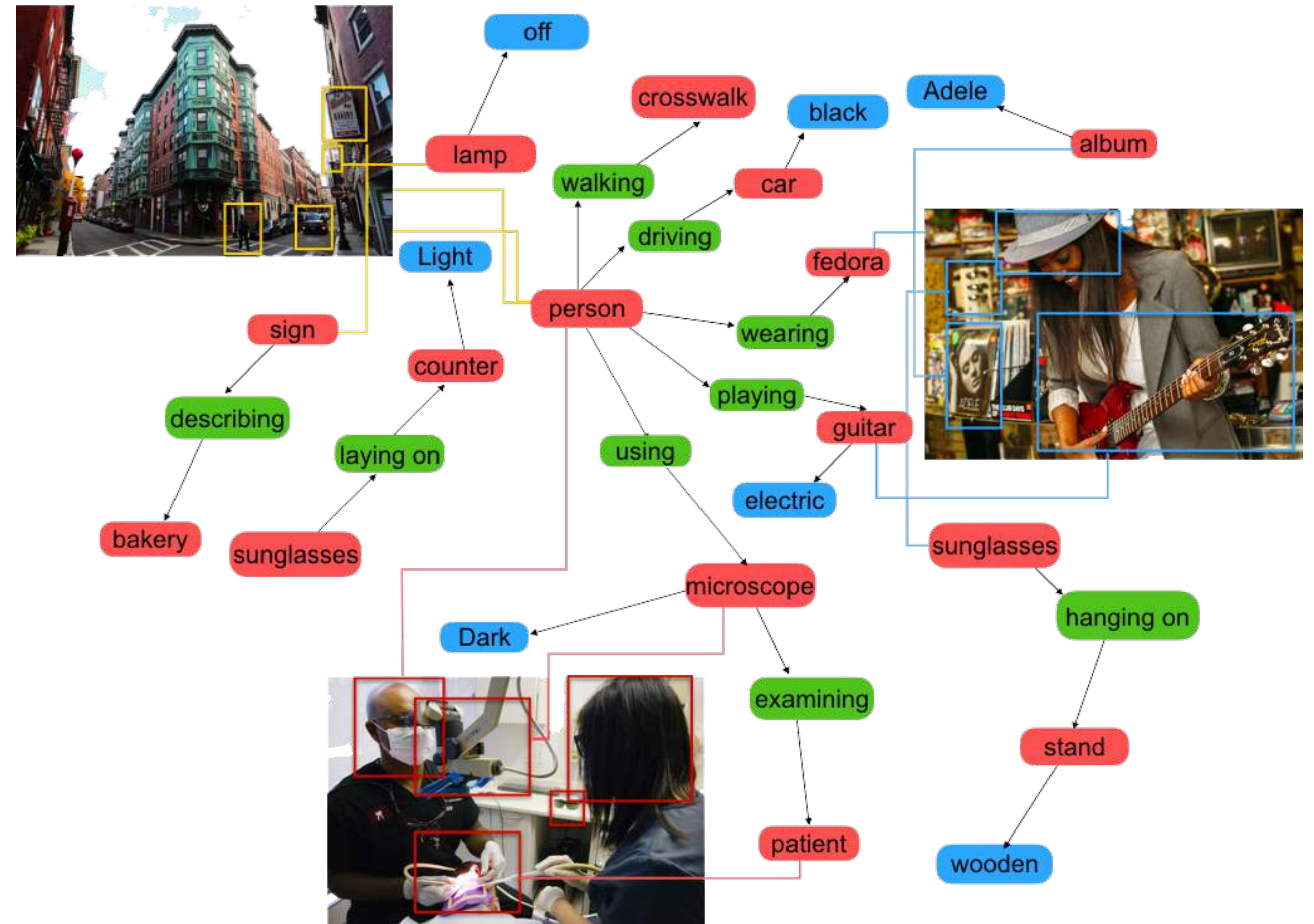Knowledge of "How-To"

Declarative Knowledge ("That-Is")

Understanding the World

Robotics

Interacting with the World

Procedural Knowledge ("How-To")

# Learning Declarative ("That-Is") Knowledge from the Web

Declarative

Knowledge

("That-Is")

Understanding
the World

Robotics

Interacting with
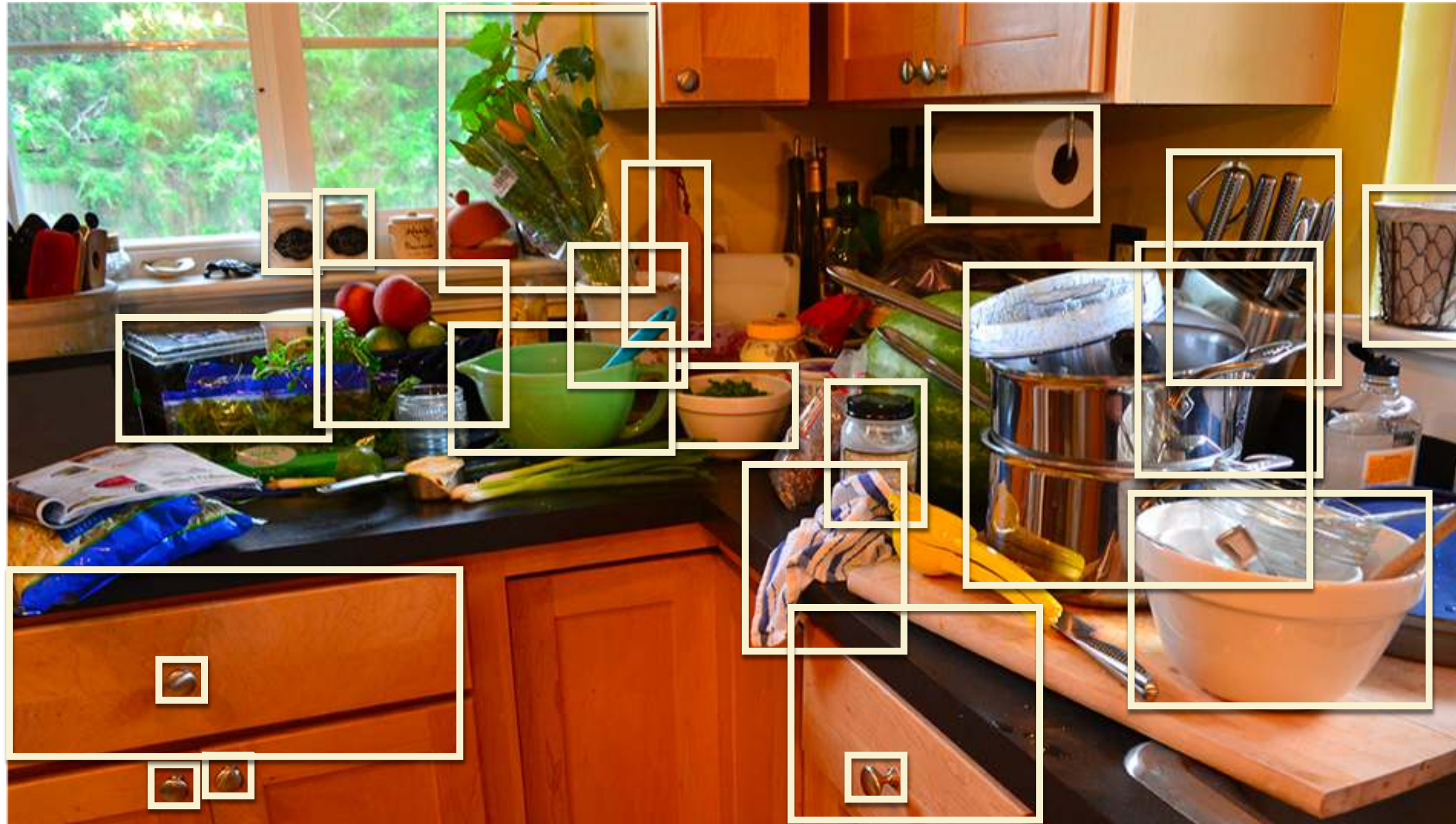the World

Procedural

Knowledge

("How-To")

Understanding the world is the cornerstone of interacting with the world.

# The Visual Genome Project

A large-scale visual knowledge base of structured image concepts

Krishna, **Zhu**, Groth, Johnson, Hata, Kravitz, Chen, Kalantidis, Li, Shamma, Bernstein, and Fei-Fei, IJCV 2017

# Visual Genome



Scene Graph: Objects + Attributes + Relationships

## Questions

1. Q: What's the color of the counter? A: Black.
2. Q: How many drawers can you see? A: Two.
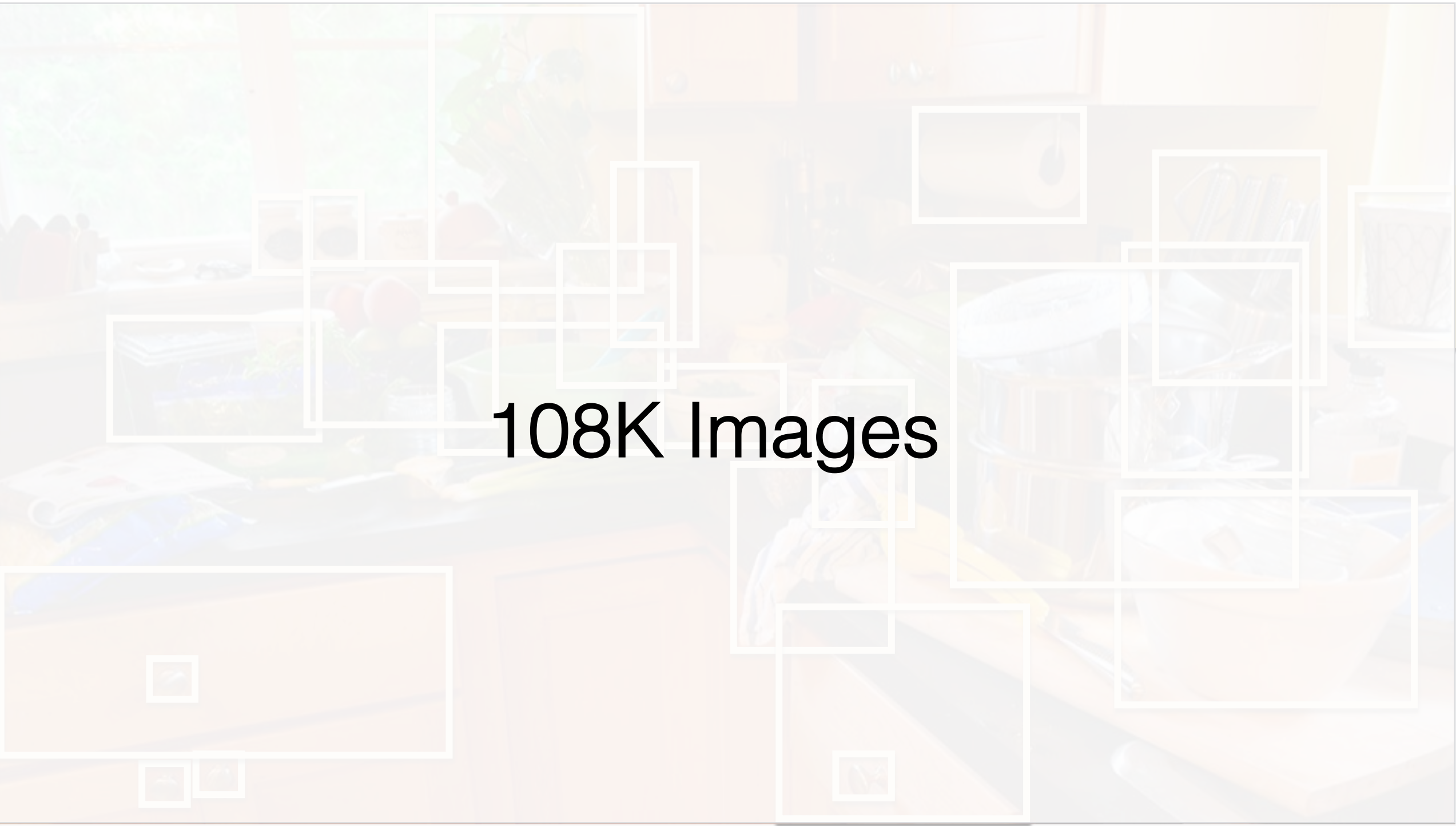3. Q: What's the material of the pots? A: Metal.
......
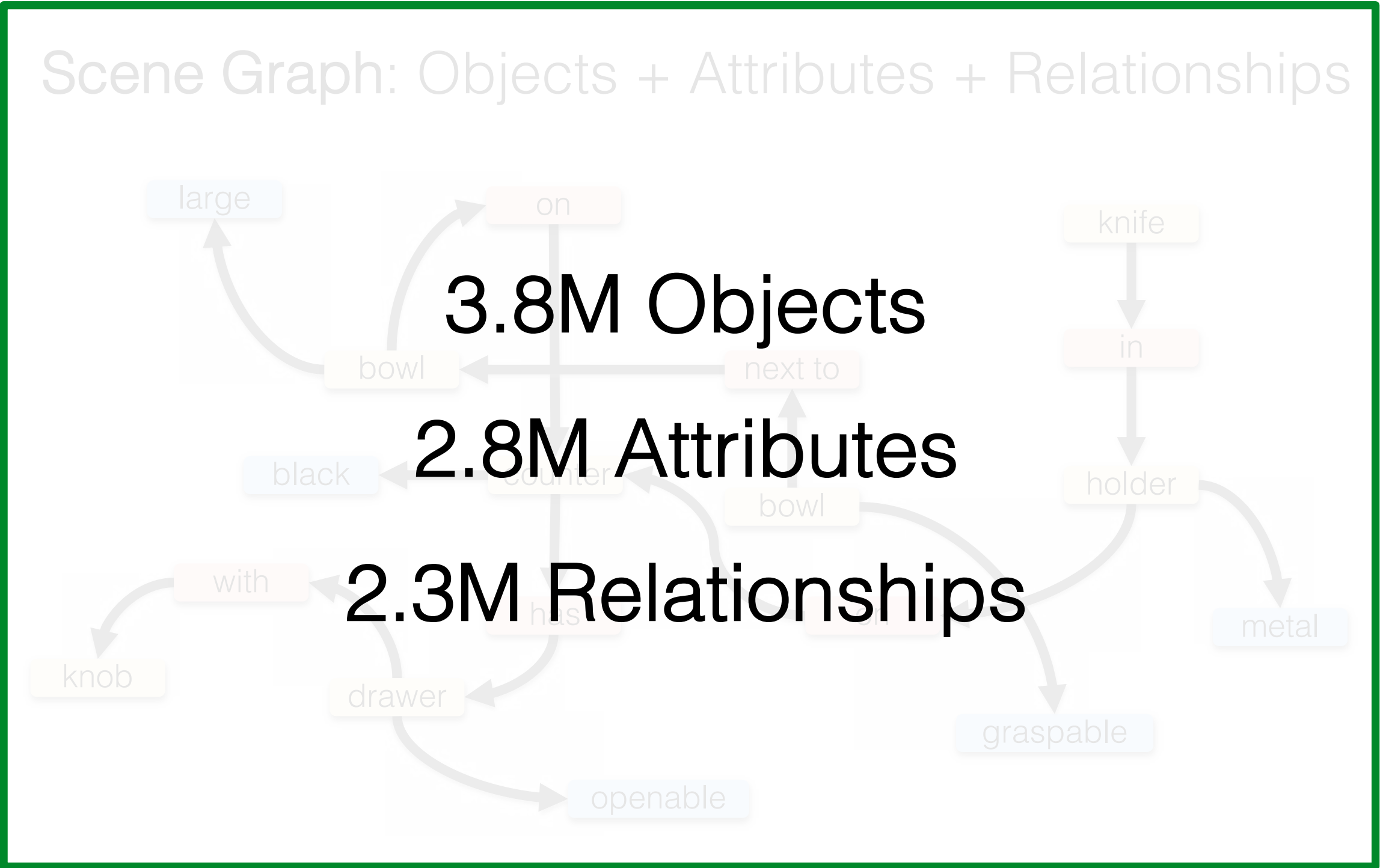
## Region Descriptions

1. There a green bowl on the black counter.
2. The cabinet door is closed.
3. Six knives are placed in the knife holder.
......

# Visual Genome



108K Images

## Scene Graph: Objects + Attributes + Relationships

3.8M Objects

2.8M Attributes

2.3M Relationships

large · on · knife · bowl · next to · in · black · counter · bowl · holder · with · metal · knob · drawer · graspable · openable

## Questions

1. Q: What's the color of the counter? A: Black.
2. Q: How many drawers can you see? A: Two.
3. Q: What's the material of the pots? A: Metal.
......

1.7M Questions

## Region Descriptions

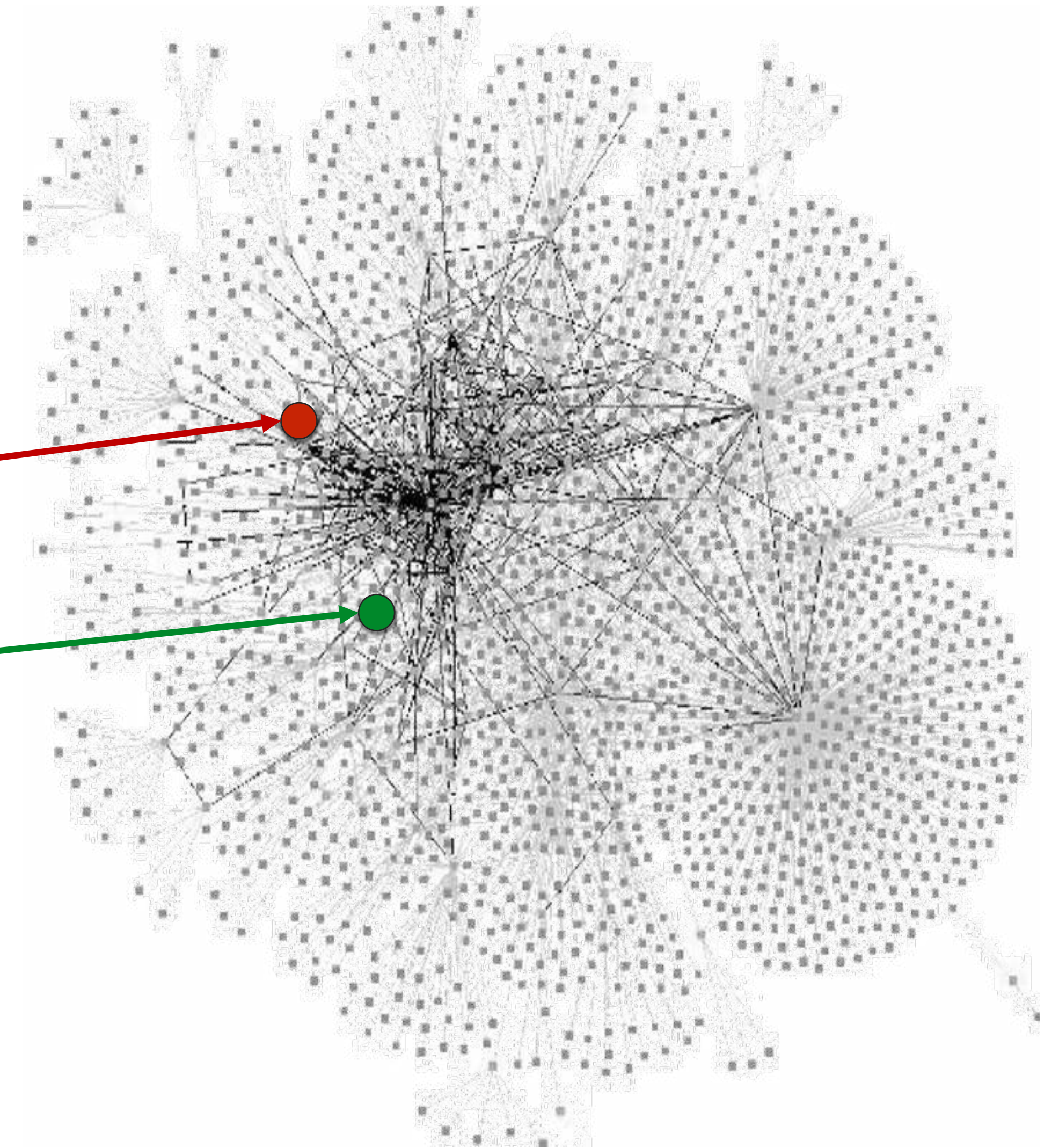1. There a green bowl on the black counter.
2. The cabinet doors
3. Six knives are placed in the knife holder.
......

5.4M Region Descriptions

# Visual Genome

An ontology of visual concepts
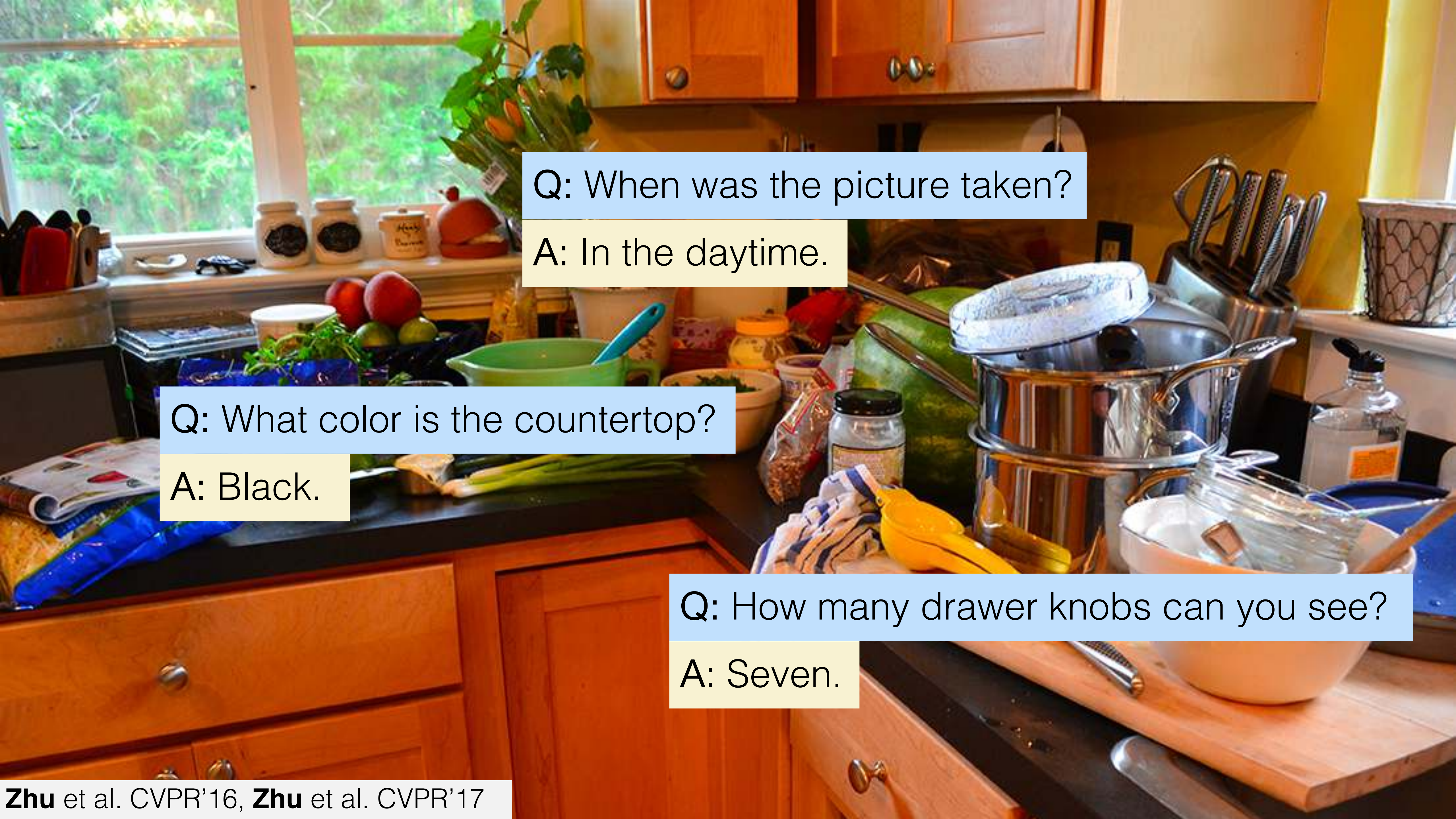
two ceramic jars

knives in a holder

green onions sitting on the counter

wooden drawer is closed

a big white bowl
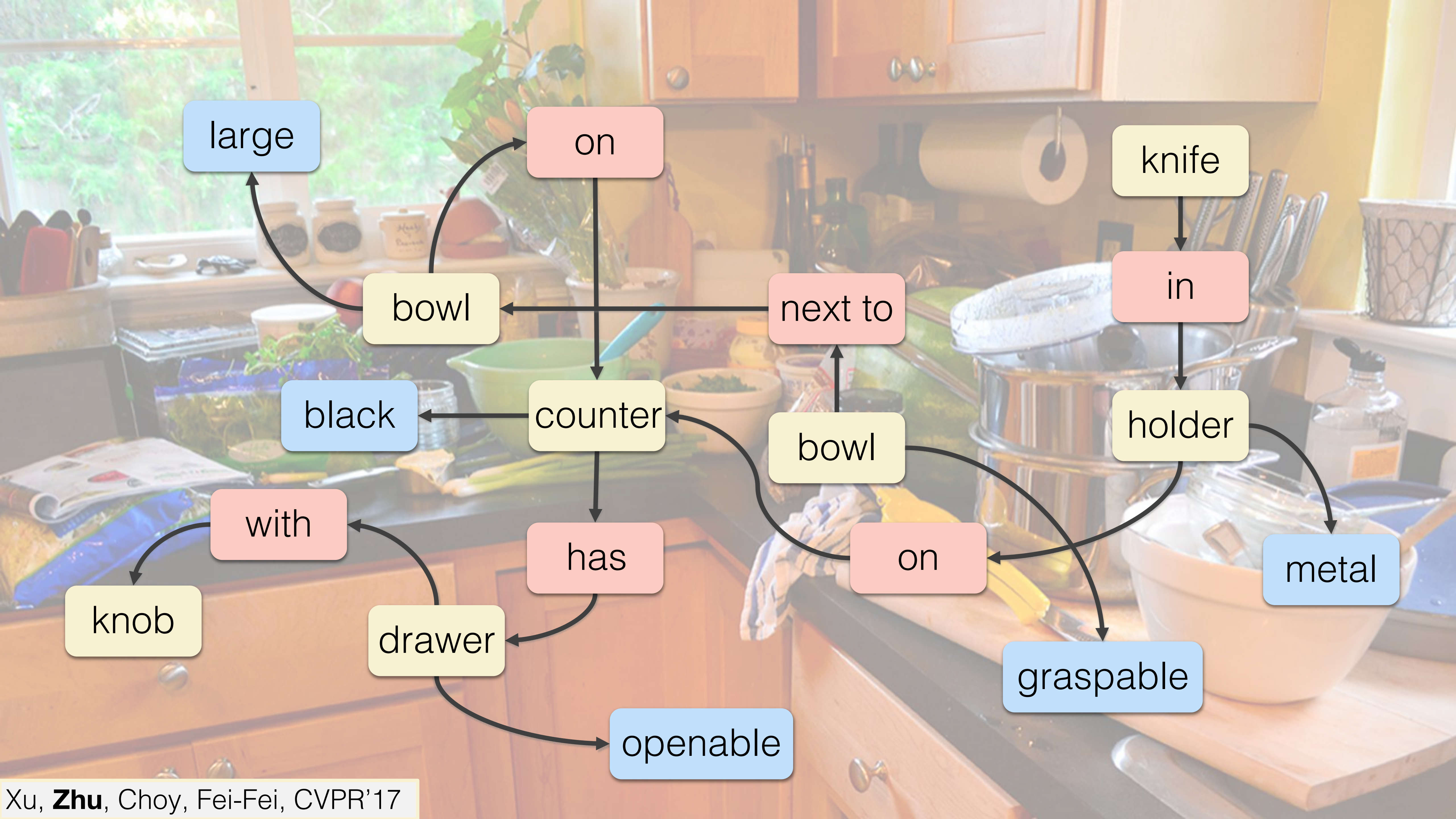
Q: When was the picture taken?
A: In the daytime.

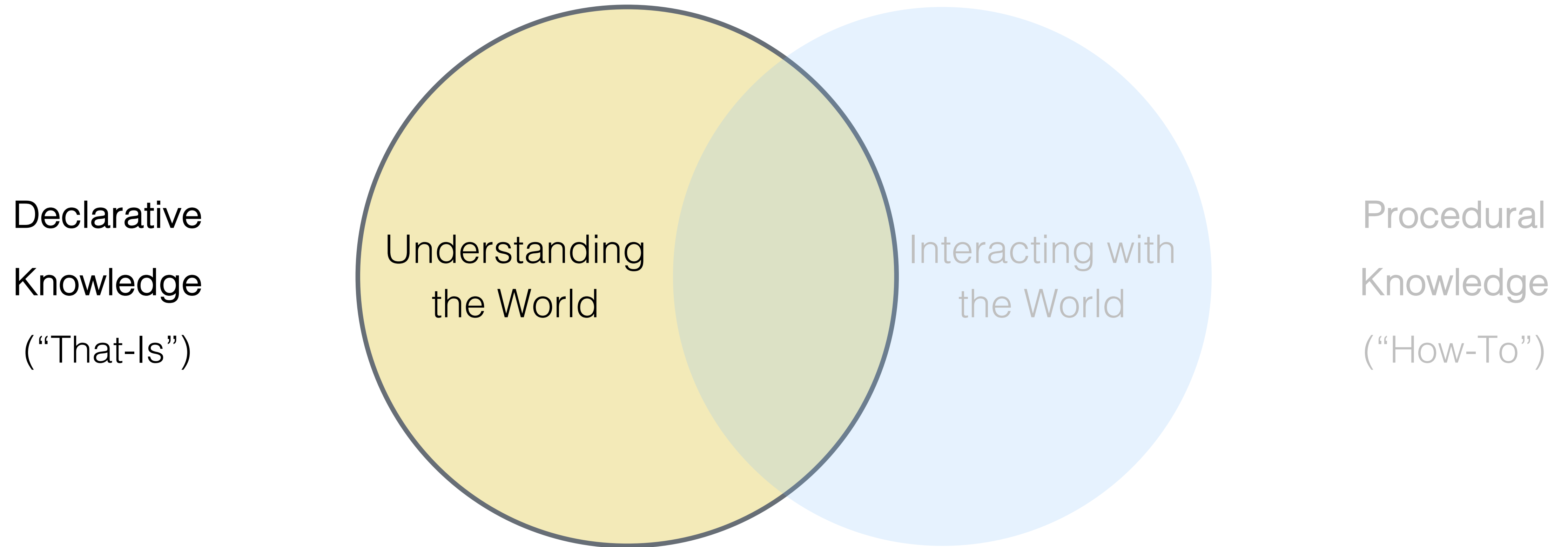Q: What color is the countertop?
A: Black.
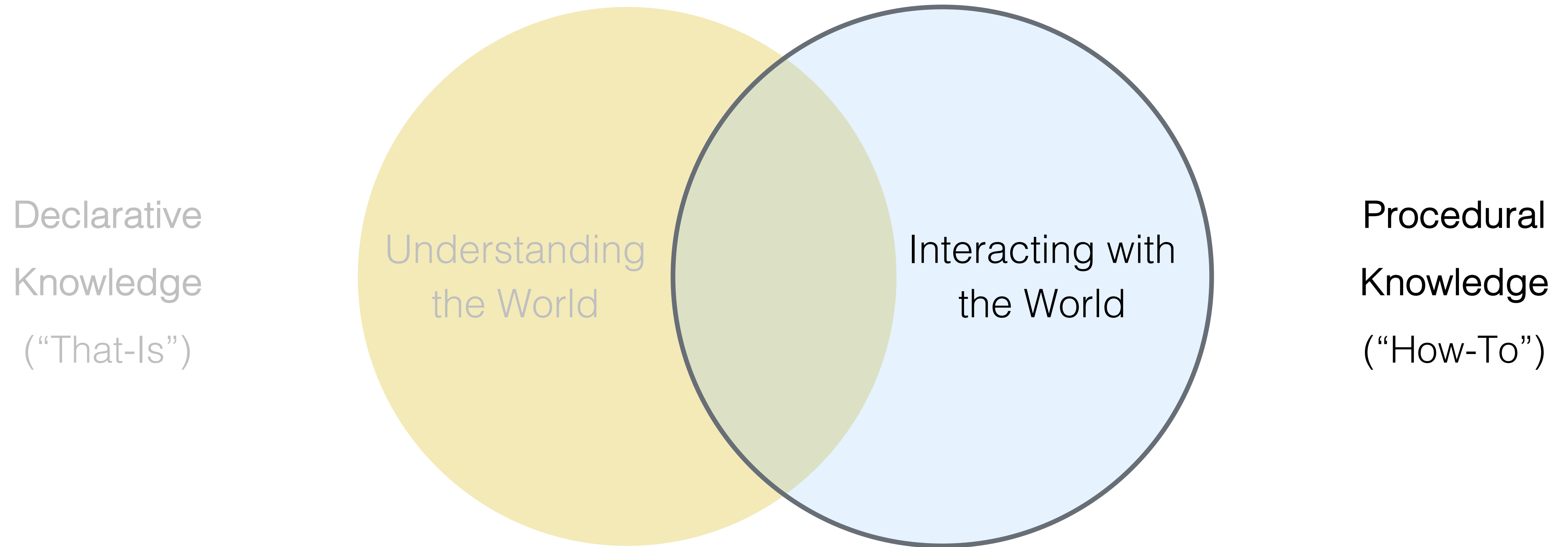
Q: How many drawer knobs can you see?
A: Seven.

**Zhu** et al. CVPR'16, **Zhu** et al. CVPR'17

Learning Procedural Knowledge needs new methodology.

It is hard to pinpoint and difficult to verbally described.

Declarative
Knowledge
("That-Is")

Understanding
the World

Interacting with
the World

Procedural
Knowledge
("How-To")

# Learning Procedural ("How-To") Knowledge from the Web

## Three Key Questions

❖ What's a good representation of procedural knowledge?

❖ How do we learn procedural knowledge from the web?

❖ How can robots take advantage of such knowledge?

# Part I: Learning from Video Demonstrations

# Part II: Learning from Crowd Teleoperation

# Part I: Learning from Video Demonstrations

# Part II: Learning from Crowd Teleoperation

# Web videos supply massive knowledge of how to solve new tasks.

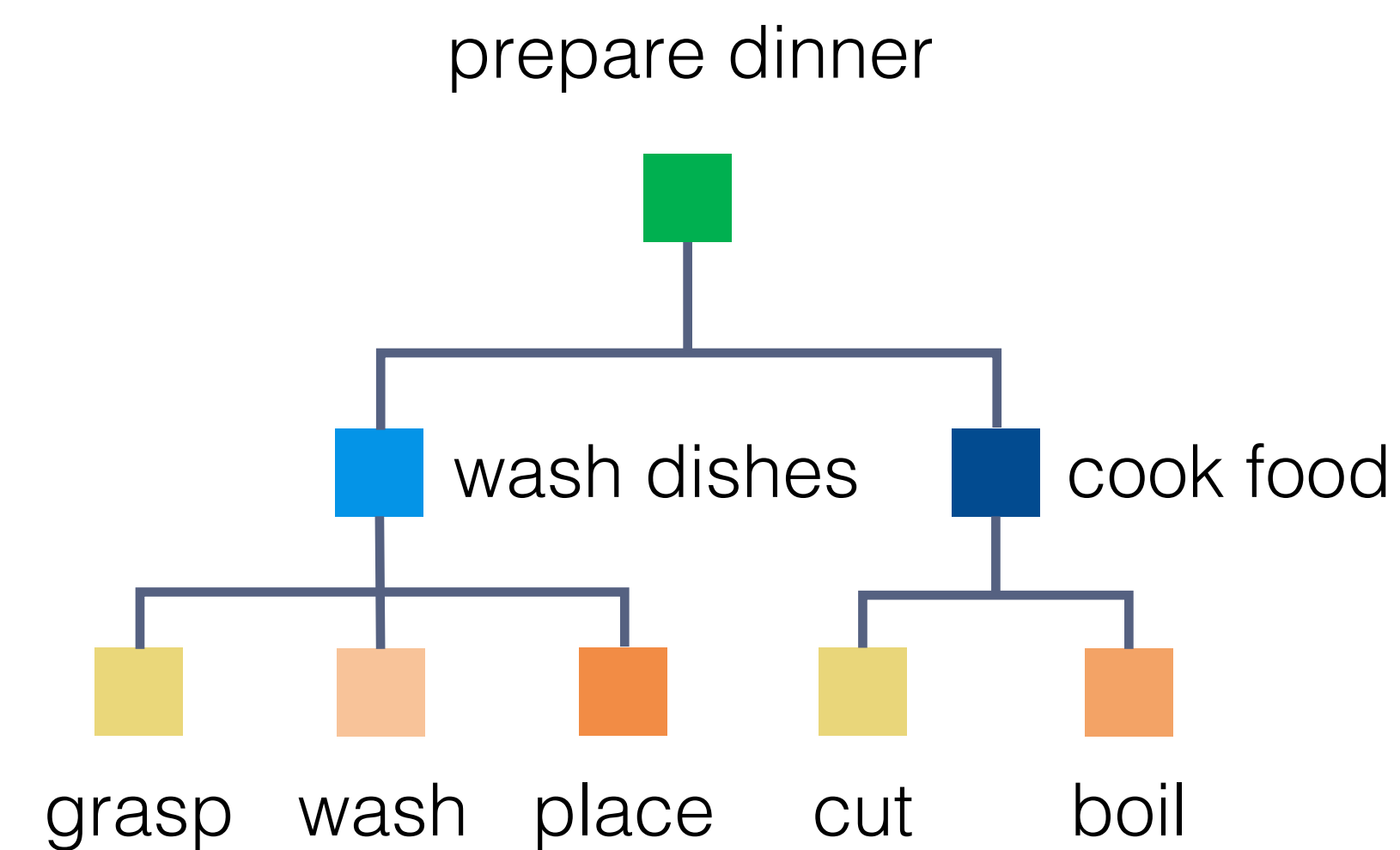# Humans learn efficiently from video demonstrations.

## Imitation of Televised Models by Infants

Andrew N. Meltzoff, *Child Development* 1988

Babies (14-24 months) can learn by imitating

demonstrations from the TV screen.





Meltzoff & Moore 1977; Meltzoff & Moore 1989, Meltzoff 1988

**Our Goal:** Learning procedural knowledge as compositional task structures from video demonstrations of a task
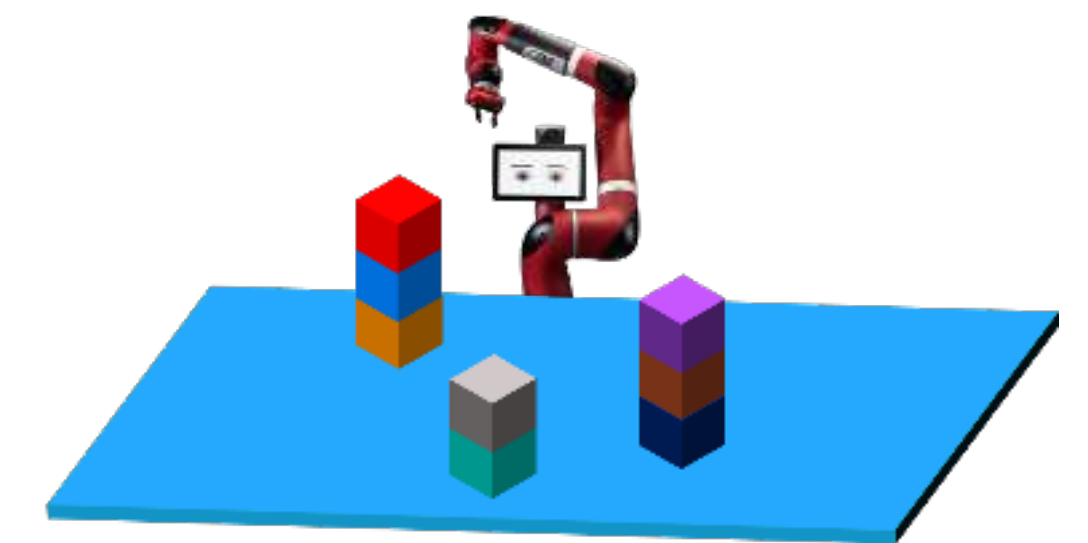


prepare dinner

# One-Shot Imitation Learning from Videos



single video
demonstration

meta-learning
model

policy for the
demonstrated task

Xu*, Nair*, **Zhu**, Gao, Garg, Fei-Fei, Savarese. *ICRA 2018*

# One-Shot Imitation Learning from Videos



supervision

meta-learning model

policy for the demonstrated task

…

a lot of training videos
(seen tasks)

Xu*, Nair*, **Zhu**, Gao, Garg, Fei-Fei, Savarese. *ICRA 2018*

# One-Shot Imitation Learning from Videos
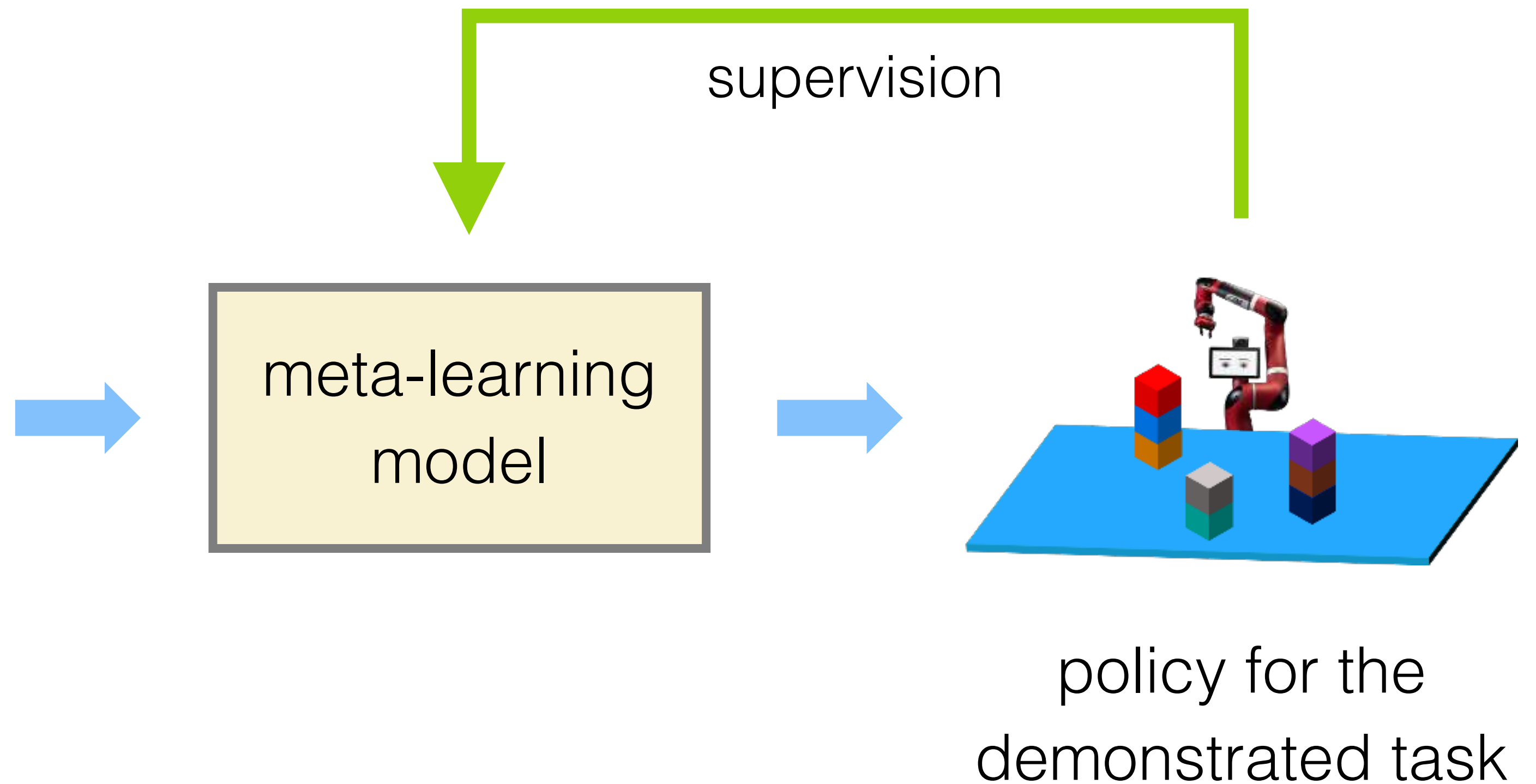


single test video
(unseen task)

meta-learning
model

policy for the
demonstrated task

# One-Shot Imitation Learning from Videos



[Duan et al. 17; Finn et al. 2017; Wang et al. 2017; Yu et al. 2018]

modeling demonstration
as a **flat sequence**

# One-Shot Imitation Learning from Videos



[Duan et al. 17; Finn et al. 2017; Wang et al. 2017; Yu et al. 2018]

modeling demonstration
as a **flat sequence**

modeling demonstration
as a **compositional structure**

# Neural Task Programming (NTP): Hierarchical Policy Learning as Neural Program Induction



Move_to (Blue)    Grip (Blue)    Move_to (Red)    Release( )

# One-Shot Imitation Learning from Videos: Neural Task Programming (NTP)



$\pi(a|x; D_1)$

$\pi(a|x; D_2)$

meta-learning model

policy

demonstration

$\pi(a|x; D_2)$

observation

end-to-end neural network (LSTM)

next program
pick(blue)

current program
pick_place(blue, green)

# One-Shot Imitation Learning from Videos: Neural Task Programming (NTP)

$\pi(a|x; D_1)$

$\pi(a|x; D_2)$

$\pi(a|x; D_1)$

$\pi(a|x; D_2)$

Training supervision

$\left\{ \left( \qquad , \qquad \right) \right\}$

video demonstration          hierarchical program trace

demonstration

observation

end-to-end neural network (LSTM)

next program
pick(blue)

current program
pick_place(blue, green)

# One-Shot Imitation Learning from Videos: Neural Task Programming (NTP)



Object Sorting

Demo

Autonomous Execution   8x

Qualitative

Better generalization with less training data than flat baselines

Quantitative

(the higher the better)

# One-Shot Imitation Learning from Videos: Neural Task Programming (NTP)



demonstration

meta-learning model

policy

Input → Black Box → Output

end-to-end
neural network
(LSTM)

compositional
model prior

Xu*, Nair*, **Zhu**, Gao, Garg, Fei-Fei, Savarese. *ICRA 2018*

# One-Shot Imitation Learning from Videos: Neural Task Graphs (NTG)

$$\pi_1(a|x;D)$$
$$\pi_2(a|x;D)$$
$$\pi_D(a|x)$$

$$\pi(a|x;D_1)$$
$$\pi(a|x;D_2)$$



demonstration

meta-learning model

policy

$\pi(a|x$

observation

$\pi(a|x$

Task Graph Generator

Neural Task Graph

Huang*, Nair*, Xu*, **Zhu**, Garg, Fei-Fei, Savarese, Niebles. *CVPR 2019*

# One-Shot Imitation Learning from Videos: Neural Task Graphs (NTG)

$$\pi_1(a|x; D)$$
$$\pi_2(a|x; D)$$
$$\pi_D(a|x)$$

$$\pi(a|x; D_1)$$
$$\pi(a|x; D_2)$$



meta-learning model

demonstration

policy

observation

Task Graph Generator

Neural Task Graph

# One-Shot Imitation Learning from Videos: Neural Task Graphs (NTG)



$$\pi_1(a|x; D)$$
$$\pi_2(a|x; D)$$
$$\pi_D(a|x)$$

$$\pi(a|x; D_1)$$
$$\pi(a|x; D_2)$$

meta-learning model

policy

demonstration

Task Graph Generator

Neural Task Graph

observation

Huang*, Nair*, Xu*, **Zhu**, Garg, Fei-Fei, Savarese, Niebles. *CVPR 2019*

# One-Shot Imitation Learning from Videos: Neural Task Graphs (NTG)

## Task Graph



place(red)

pick(orange)

pick(green)

pick(red)

place(green)

| Nodes | States | **infinite** |
|-------|--------|----------|
| Edges | Actions | |

## Conjugate Task Graph



place(green)  pick(green)

valid states

pick(orange)

pick(red)  place(red)

| Nodes | Actions | **finite** |
|-------|---------|----------|
| Edges | States (Preconditions) | |

Huang*, Nair*, Xu*, **Zhu**, Garg, Fei-Fei, Savarese, Niebles. *CVPR 2019*

# One-Shot Imitation Learning from Videos: Neural Task Graphs (NTG)



current observation

node localizer

edge classifier

selected **node**

place(green)

selected **edge**

pick(red)

pick(green)

pick(orange)

place(red)

pick(red)

next action

Huang*, Nair*, Xu*, **Zhu**, Garg, Fei-Fei, Savarese, Niebles. *CVPR 2019*

# One-Shot Imitation Learning from Videos: Neural Task Graphs (NTG)

$$\pi_1(a|x;D)$$
$$\pi_2(a|x;D)$$
$$\pi_D(a|x)$$

$$\pi(a|x;D_1)$$
$$\pi(a|x;D_2)$$



Training supervision

$$\Big\{\Big(\ \ \text{video demonstration} \ \ \ \ \ \text{uence}\ \ \Big)\Big\}$$

$\pi(a|x$

demonstration

policy

$\pi(a|x$

observation

$\pi(a|x$

Task Graph Generator

Neural Task Graph

Huang*, Nair*, Xu*, **Zhu**, Garg, Fei-Fei, Savarese, Niebles. *CVPR 2019*

# One-Shot Imitation Learning from Videos: Neural Task Graphs (NTG)



Recovery from Intermediate Failures

Autonomous Execution                                                    20x

Qualitative



Weaker supervision, less training data, and better generalization

Quantitative

(the higher the better)

Huang*, Nair*, Xu*, **Zhu**, Garg, Fei-Fei, Savarese, Niebles. *CVPR 2019*
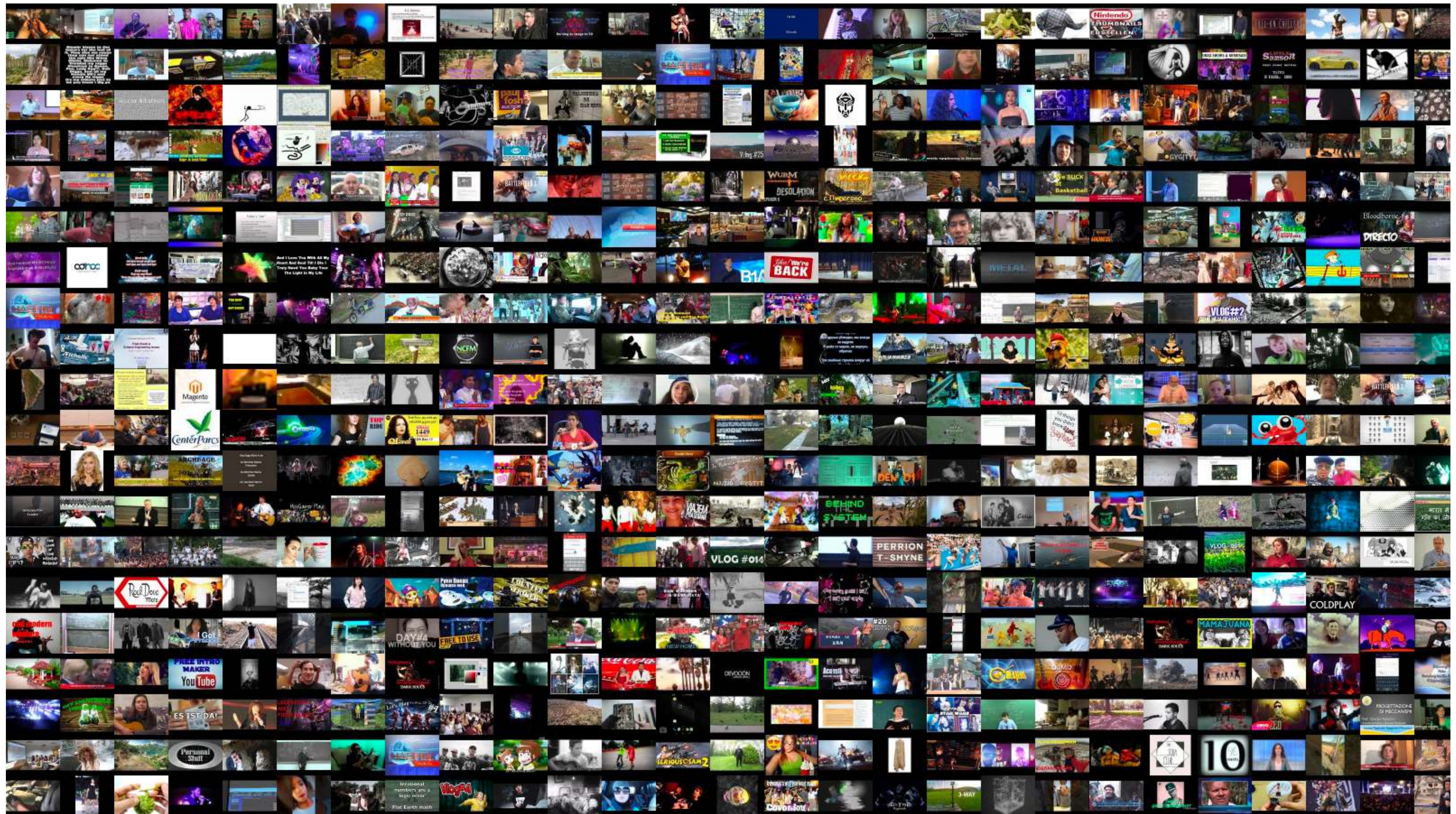
# One-Shot Imitation Learning from Videos: Neural Task Graphs (NTG)



Applying NTG to the real-world surgical video dataset JIGSAWS

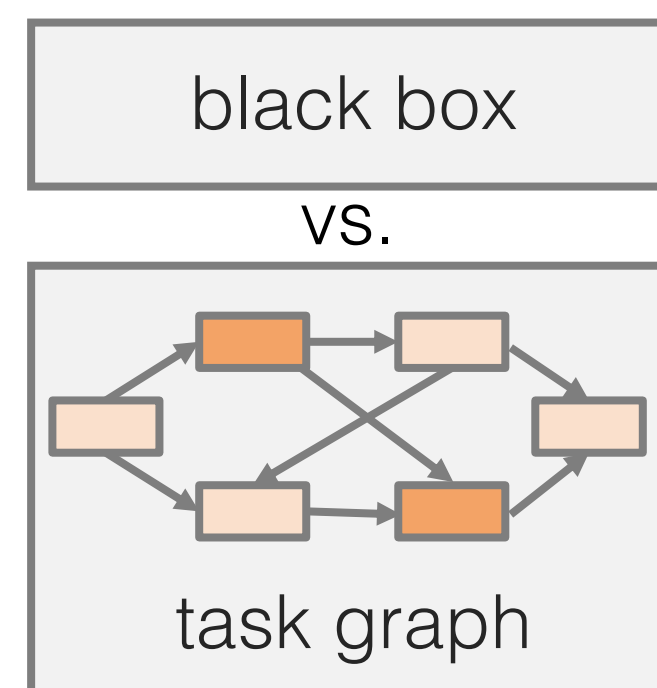Huang*, Nair*, Xu*, **Zhu**, Garg, Fei-Fei, Savarese, Niebles. *CVPR 2019*

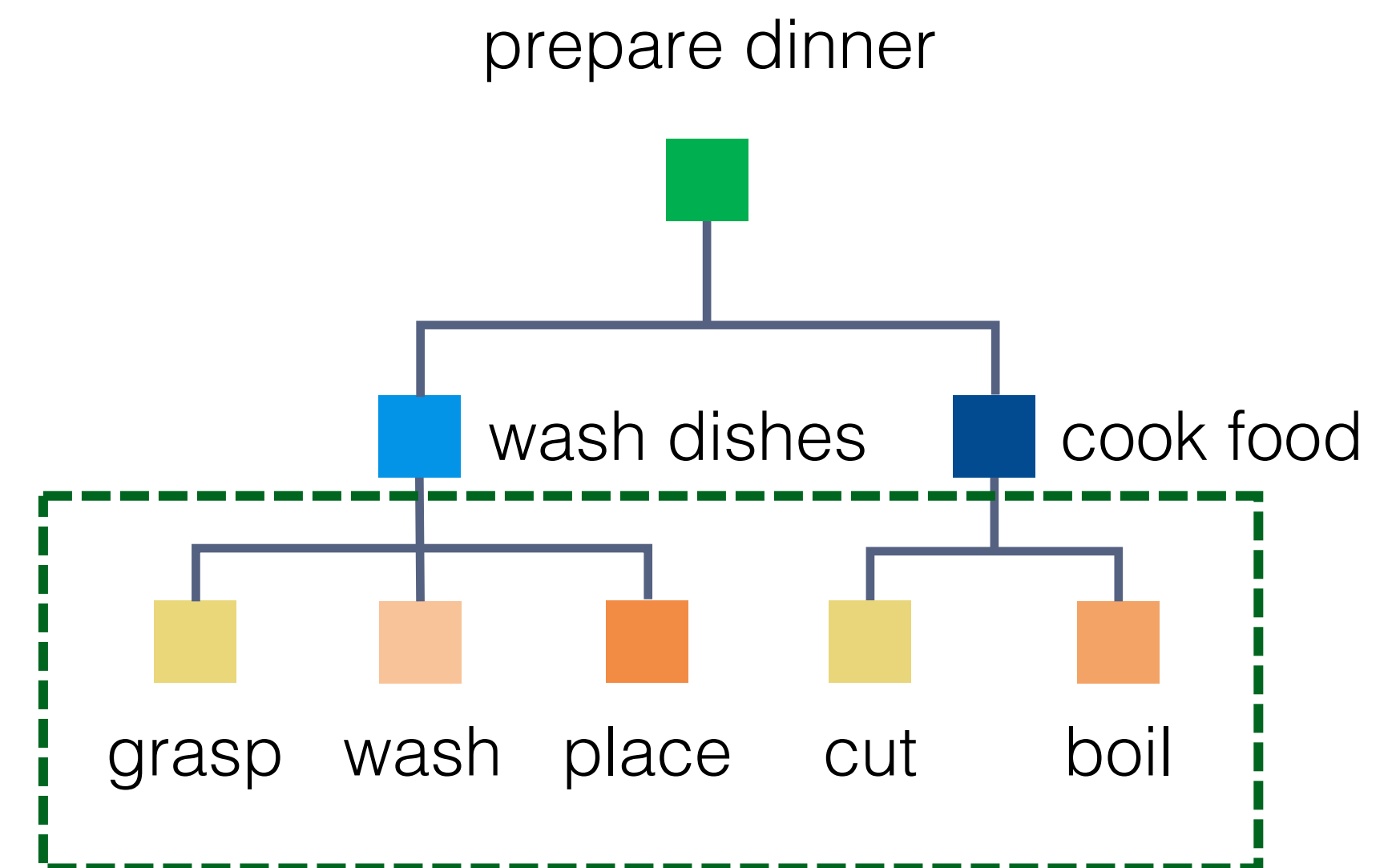**Next Goal:** Learning task knowledge from web videos

Extracting how-to knowledge about the **compositional task structure** of complex tasks from **video demonstrations**



**Meta-learning** models with compositional priors generalize better than black-box models

NTP and NTG learn how-to knowledge in the form of compositional task structures while motor skills are abstracted away.



prepare dinner

prepare dinner

wash dishes     cook food

grasp  wash  place    cut     boil

modeled as pre-defined "API calls"

NTP and NTG learn how-to knowledge in the form of compositional task structures while motor skills are abstracted away.



Manually defining motor skills is intractable.
We need to learn from data.

How can we collect data for learning motor skills from the web?
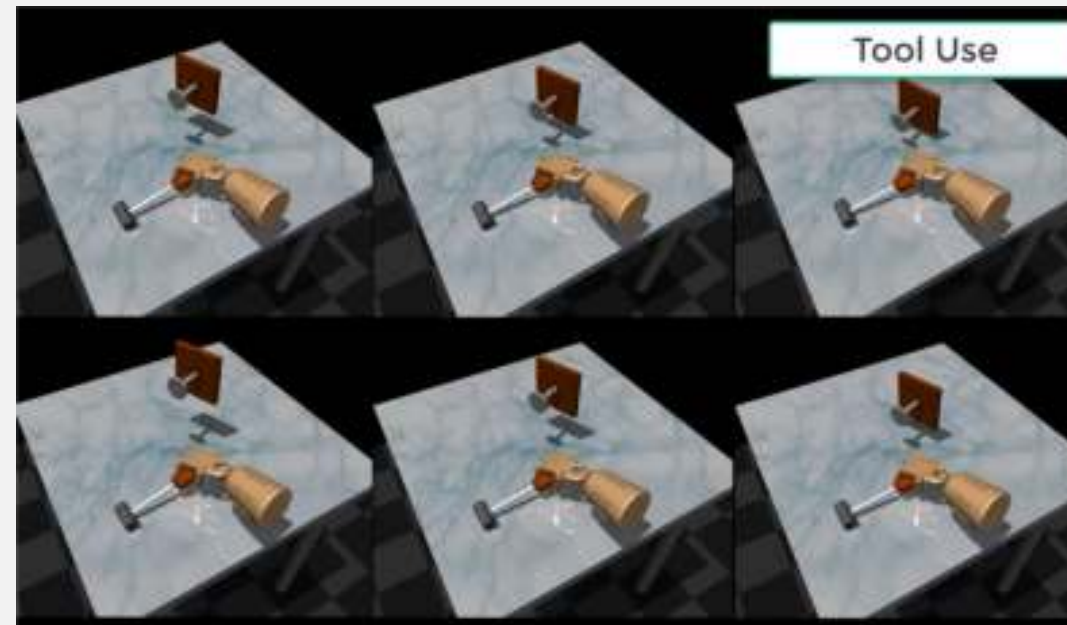
Part I: Learning from Video Demonstrations

Part II: Learning from Crowd Teleoperation

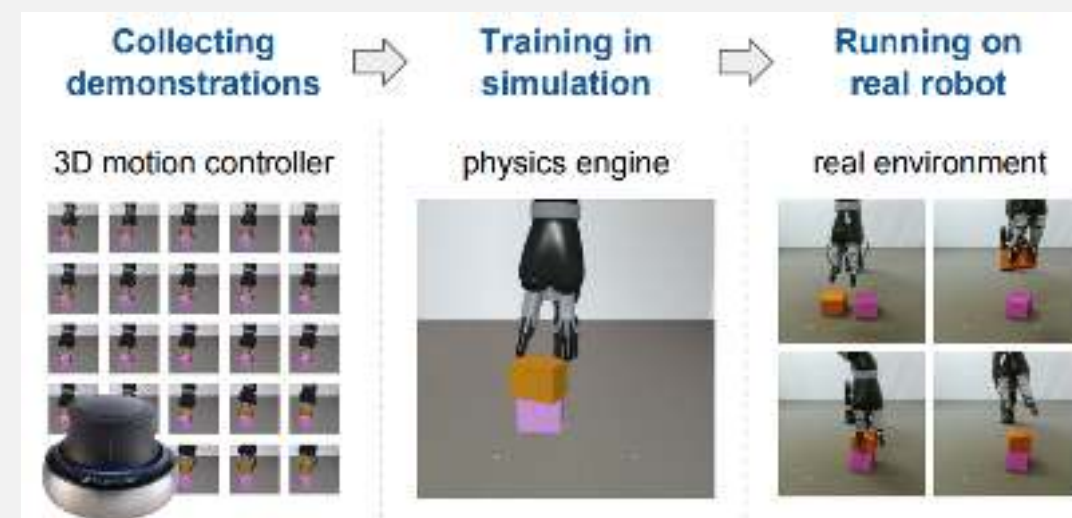# Data is critical for learning robot motor skills.

## Imitation Learning

Rajeswaran et al. 2018: 25 demos

Finn et al. 2017: 30 demos

Zhu et al. 2018: 30 demos

Vecerik et al. 2017: 100 demos

## Reinforcement & Self-Supervised Learning

Levine et al. 2016

Kalashnikov et al. 2018

Pinto et al. 2016

Fang et al. 2018

Large demonstration datasets is hard to collect.
Humans need to demonstrate not label.

Data can be low quality due to lack of expert.

Data is critical for learning robot skills.

How to scale up high-quality human supervision for robotics?

Provide a natural way for **anyone** to provide demonstrations

# Web-based Crowd Teleoperation with RoboTurk

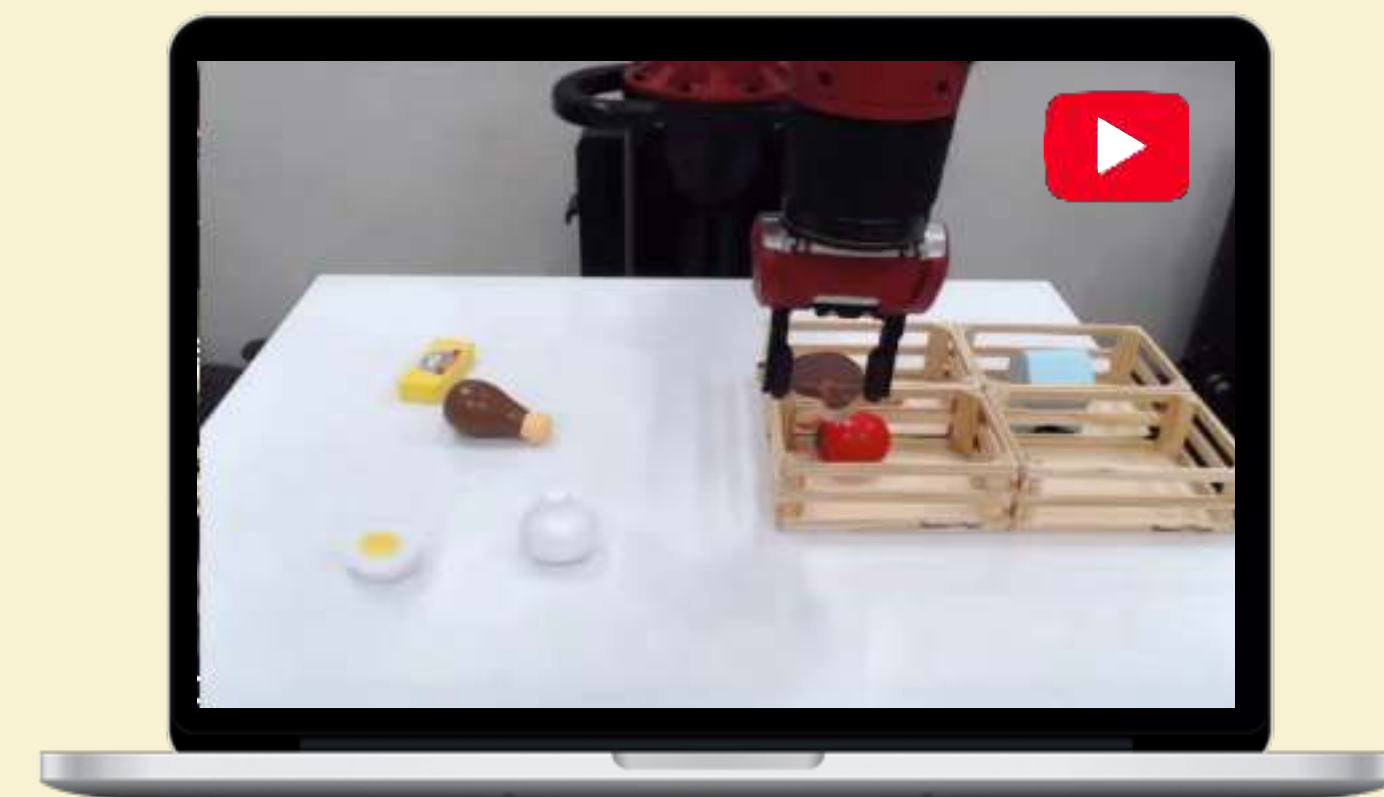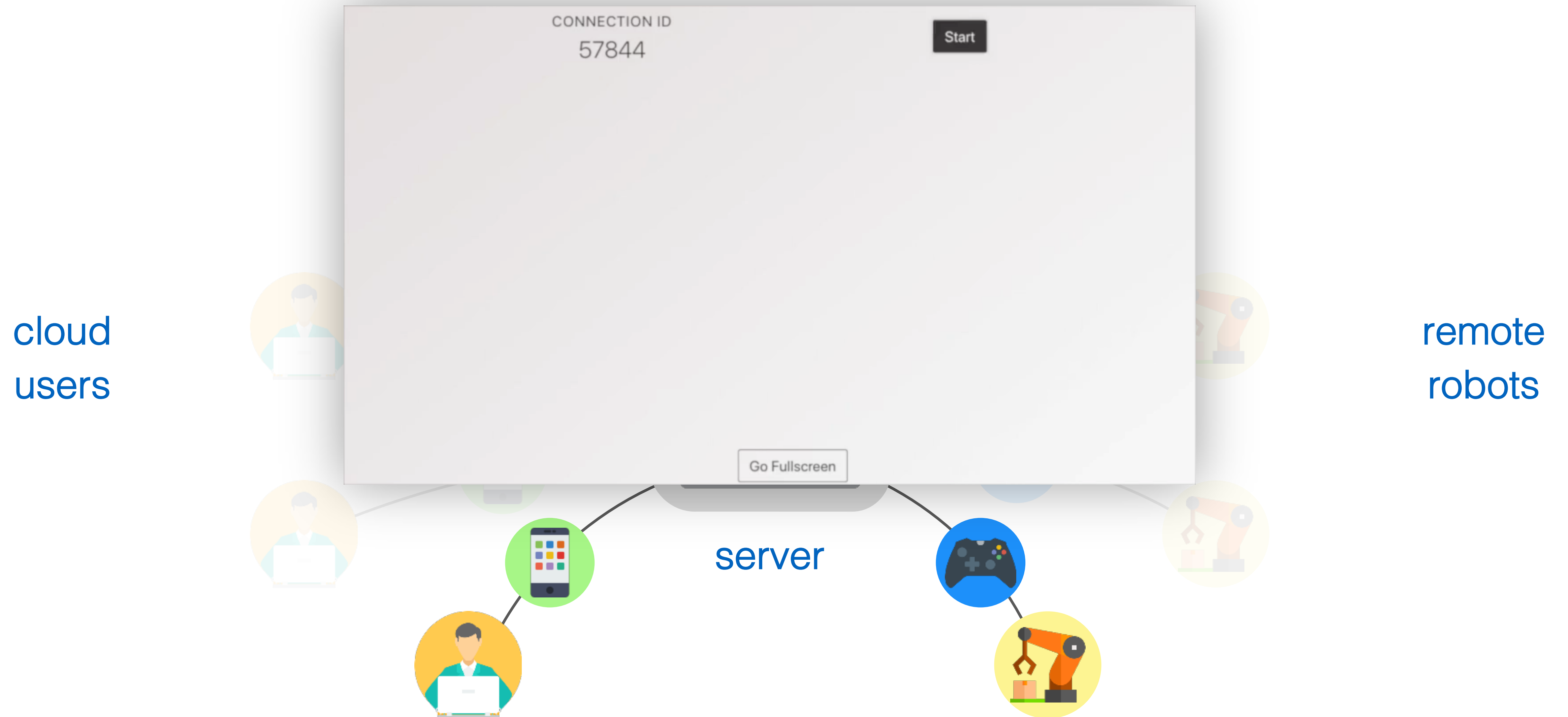## RoboTurk: Crowdsourcing Platform for Large-Scale Demonstration Collection



RoboTurk in action



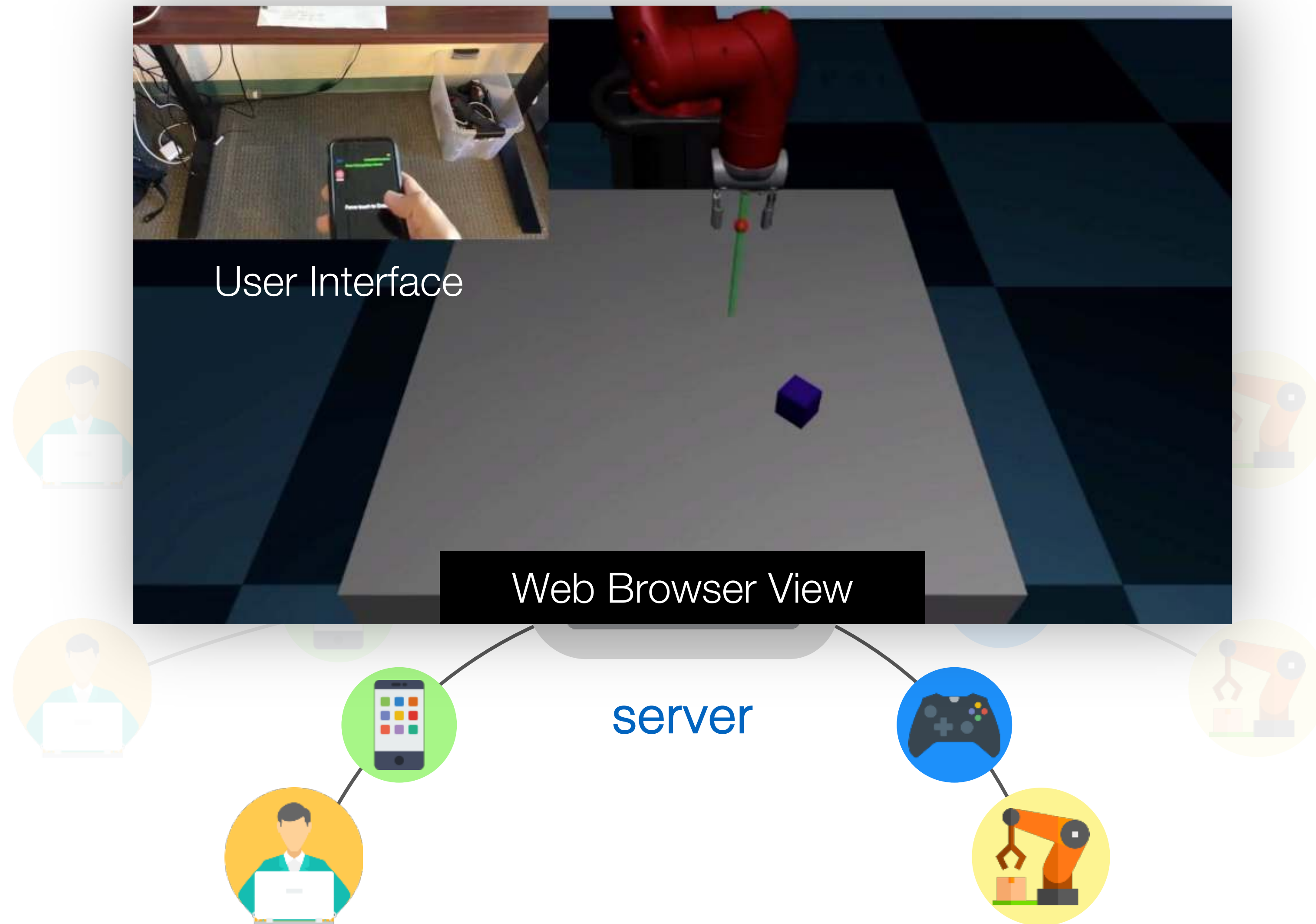6-DoF controller    +    real-time streaming from remote robot

# Web-based Crowd Teleoperation with RoboTurk

**CONNECTION ID**
57844

Start

Go Fullscreen

cloud
users

remote
robots

server

Mandlekar, **Zhu**, Garg, Booher, Spero, Tung, Gao, Emmons, Gupta, Orbay, Savarese, Fei-Fei, CoRL 2018

# Web-based Crowd Teleoperation with RoboTurk



User Interface

Web Browser View

cloud
users

remote
robots

server

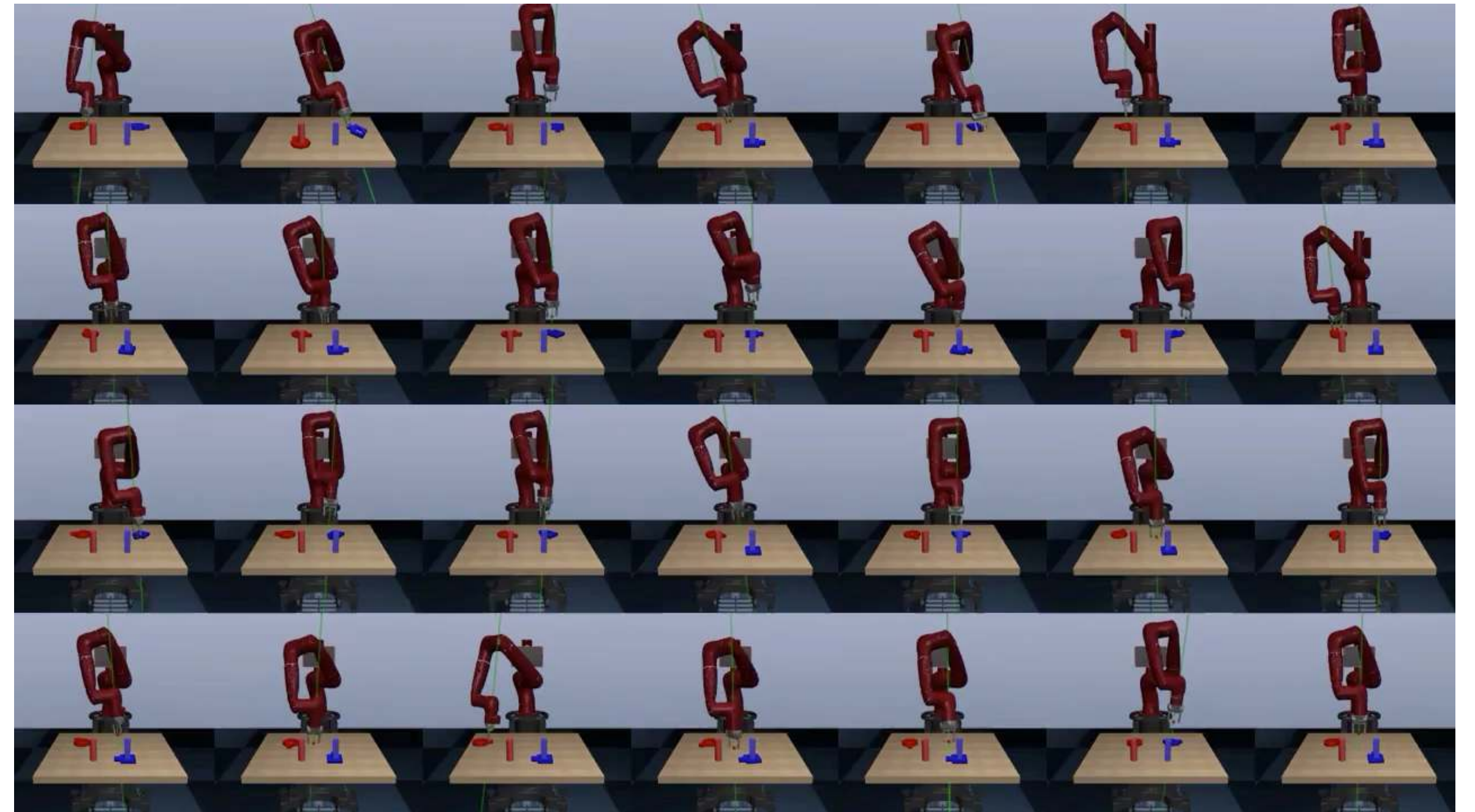Mandlekar, **Zhu**, Garg, Booher, Spero, Tung, Gao, Emmons, Gupta, Orbay, Savarese, Fei-Fei, CoRL 2018

# Web-based Crowd Teleoperation with RoboTurk

## RoboTurk Pilot Dataset

**137.5 hours** of demonstrations

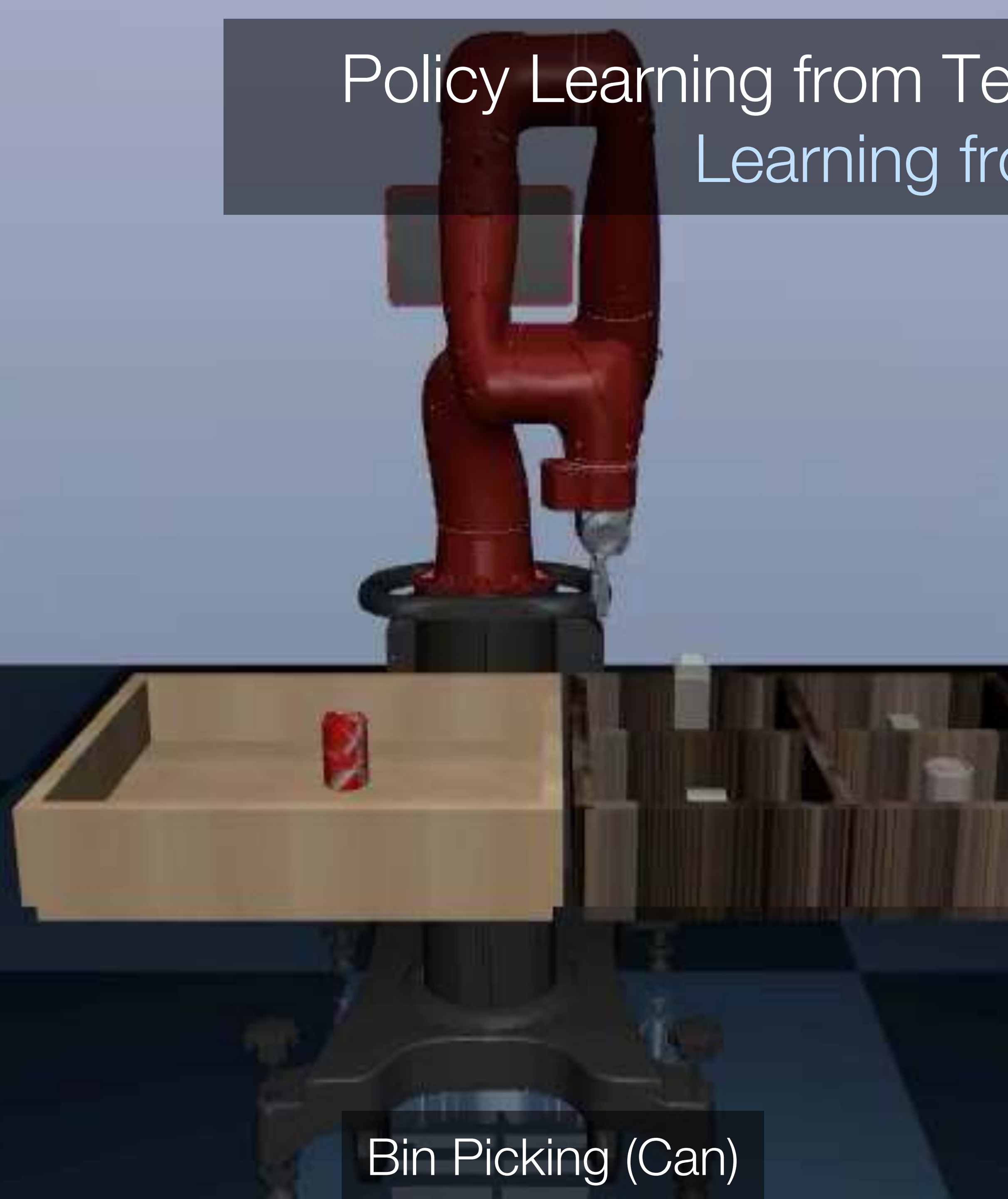**22 hours** of total platform usage
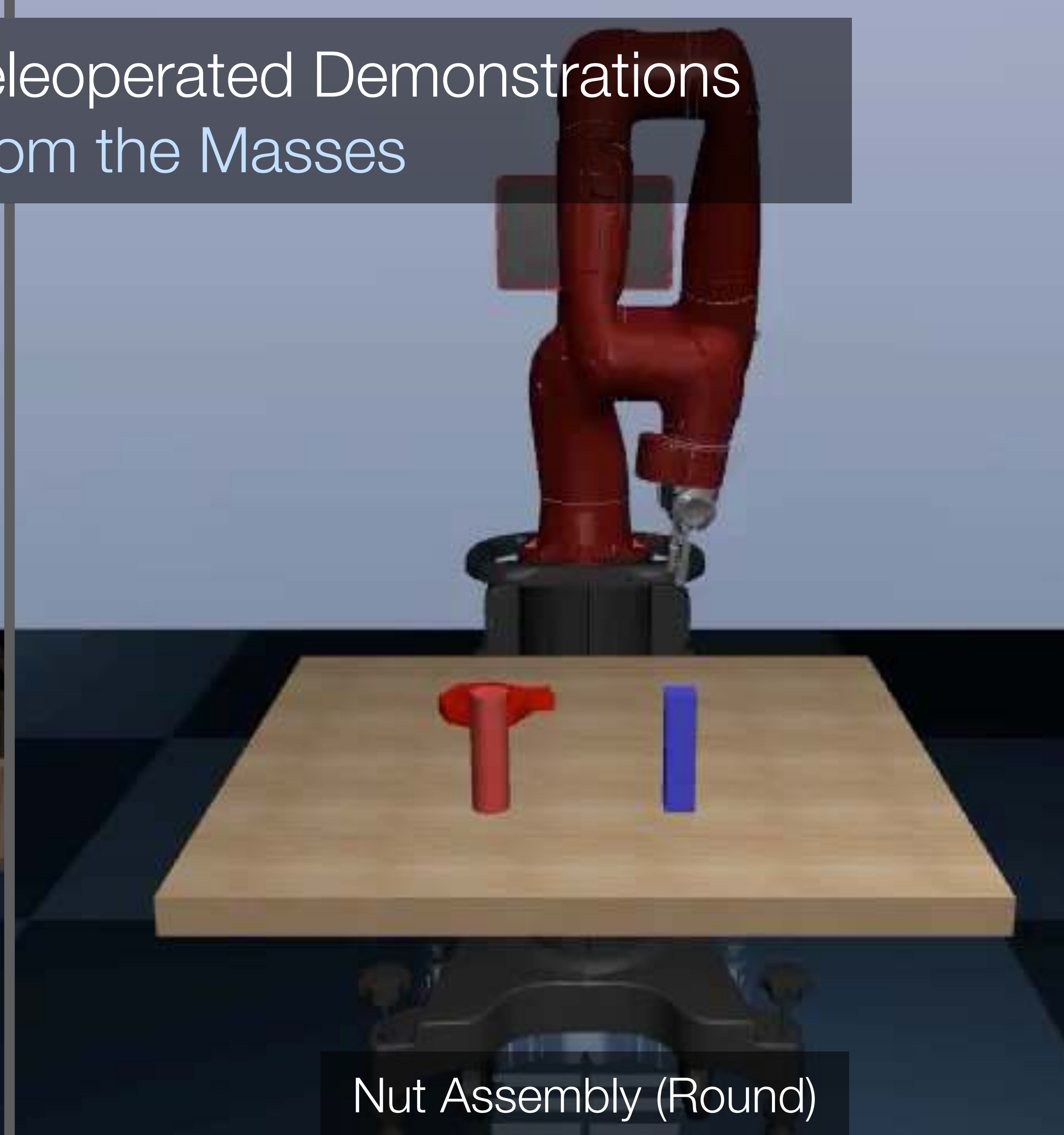
**2218** successful demonstrations



teleoperated demonstrations

Policy Learning from Teleoperated Demonstrations
Learning from the Masses

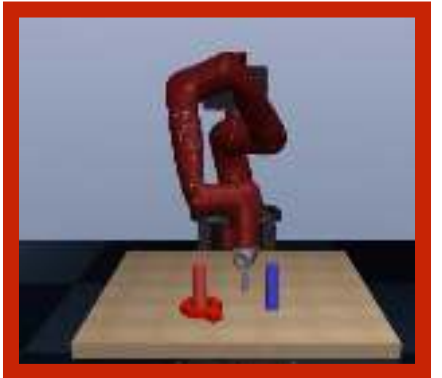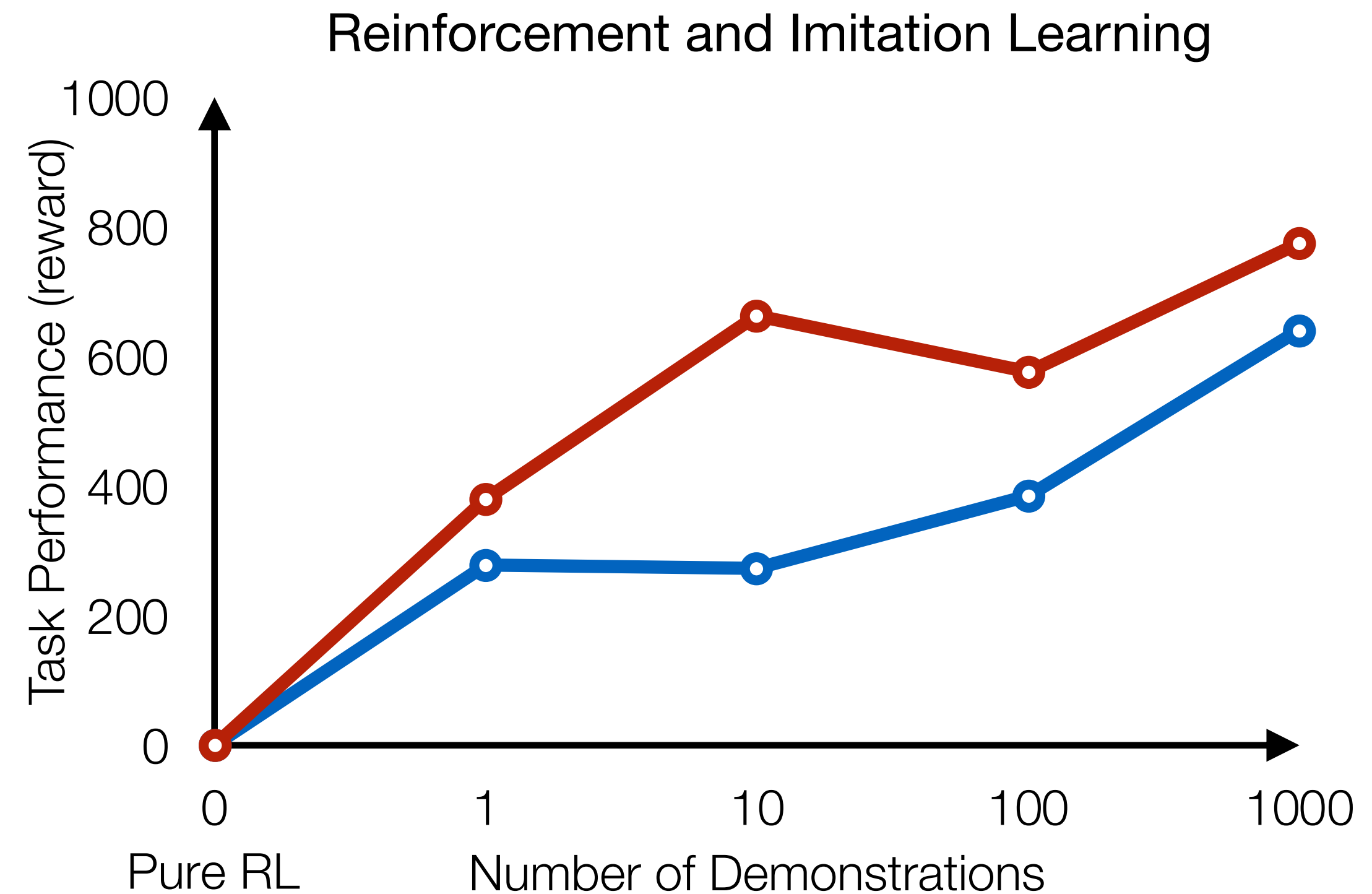Bin Picking (Can)

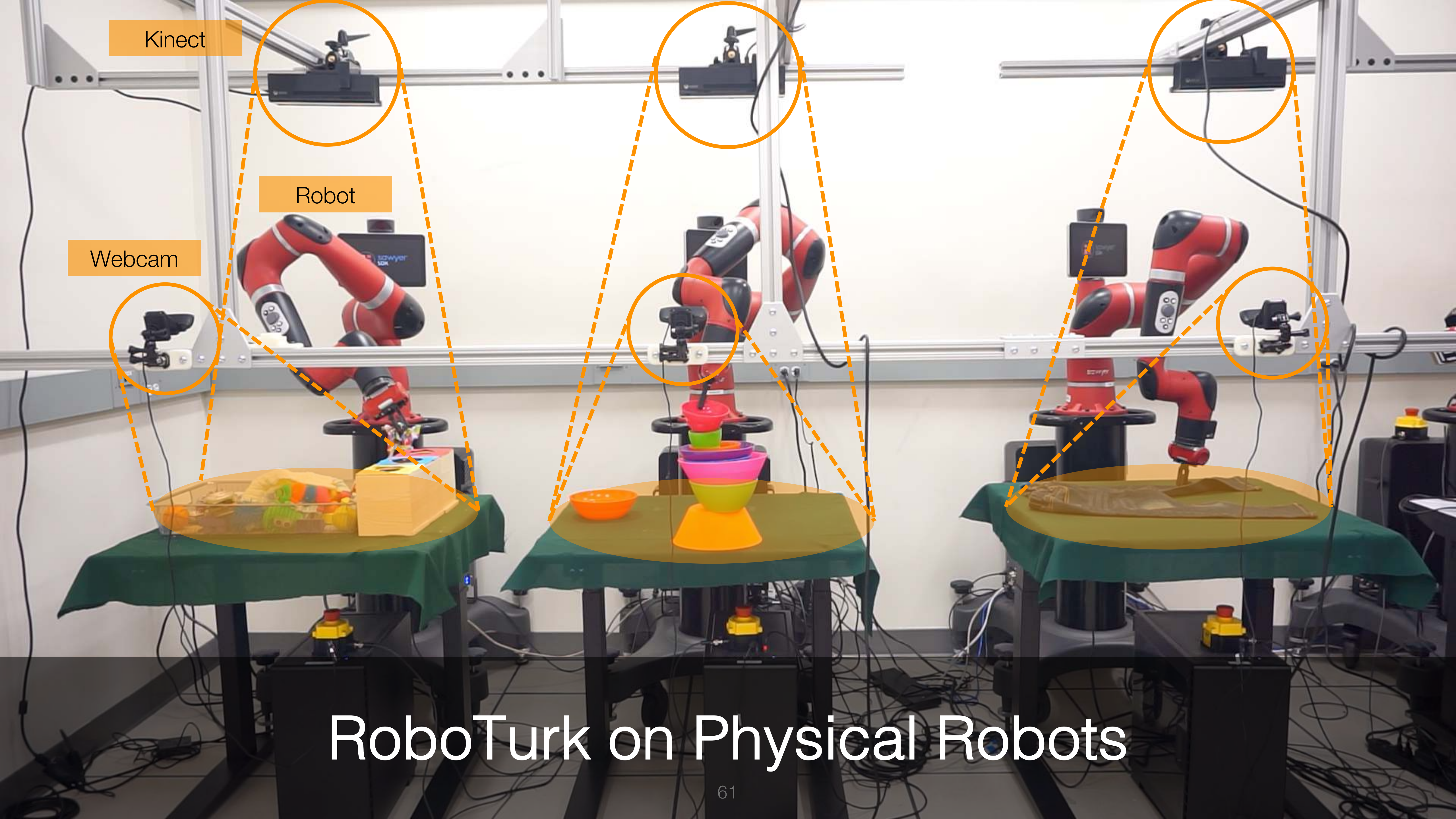Nut Assembly (Round)

# Reinforcement and Imitation Learning: Data

## RoboTurk Pilot Dataset

**137.5 hours** of demonstrations

**22 hours** of total platform usage

**2218** successful demonstrations



Reinforcement and Imitation Learning

Task Performance (reward) vs Number of Demonstrations

Zhu*, Fan*, Zhu, Liu, Zeng, Gupta, Creus-Costa, Savarese, Fei-Fei, CoRL 2018

Mandlekar, **Zhu**, Garg, Booher, Spero, Tung, Gao, Emmons, Gupta, Orbay, Savarese, Fei-Fei, CoRL 2018

Kinect

Robot

Webcam
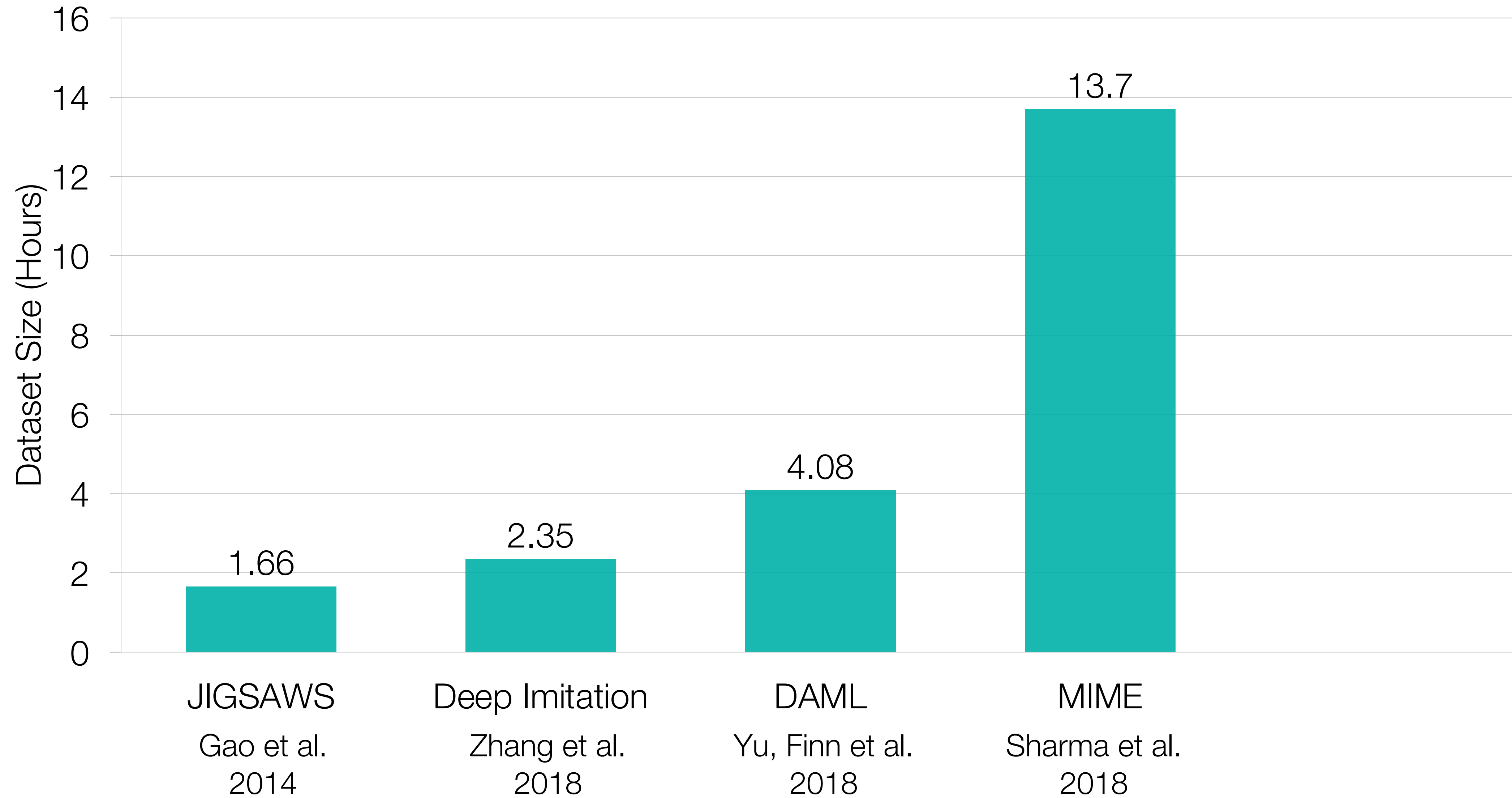
RoboTurk on Physical Robots

61

Real-time

Scalable Data Collection

# Dataset Size Comparison



Mandlekar, Booher, Spero, Tung, Gupta, **Zhu**, Garg, Savarese, Fei-Fei, IROS 2019

# Dataset Size Comparison



Mandlekar, Booher, Spero, Tung, Gupta, **Zhu**, Garg, Savarese, Fei-Fei, IROS 2019
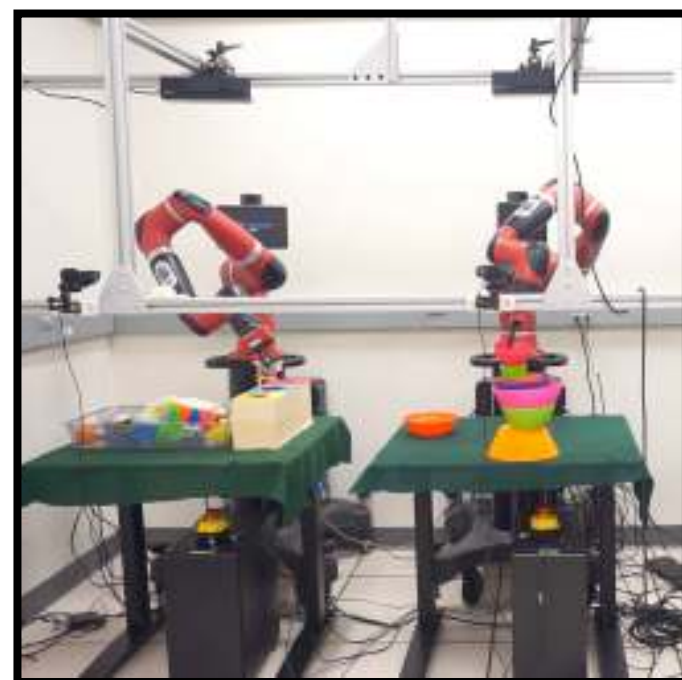
RoboTurk for

everyone, everywhere

RoboTurk scales up demonstration collection with **teleoperated crowdsourcing** from web users



Large-scale **crowdsourced data** enables us to train more effective **motor skill learning** algorithms.

# Learn More about RoboTurk?

Come to our IROS Presentation

# RoboTurk: Human Reasoning and Dexterity for Large-Scale Dataset Creation

Tuesday 15:45-16:00, Award Session II: Paper TuBT4.5

# Part I: Learning from Web Videos

Extracting compositional task structures from video data

# Part II: Learning from Crowd Teleoperation

Crowdsourcing teleoperated demonstrations for skill learning

# Conclusions

❖ What's a good representation of procedural knowledge?

High-level task structures & low-level motor skills

❖ How do we learn procedural knowledge from the web?

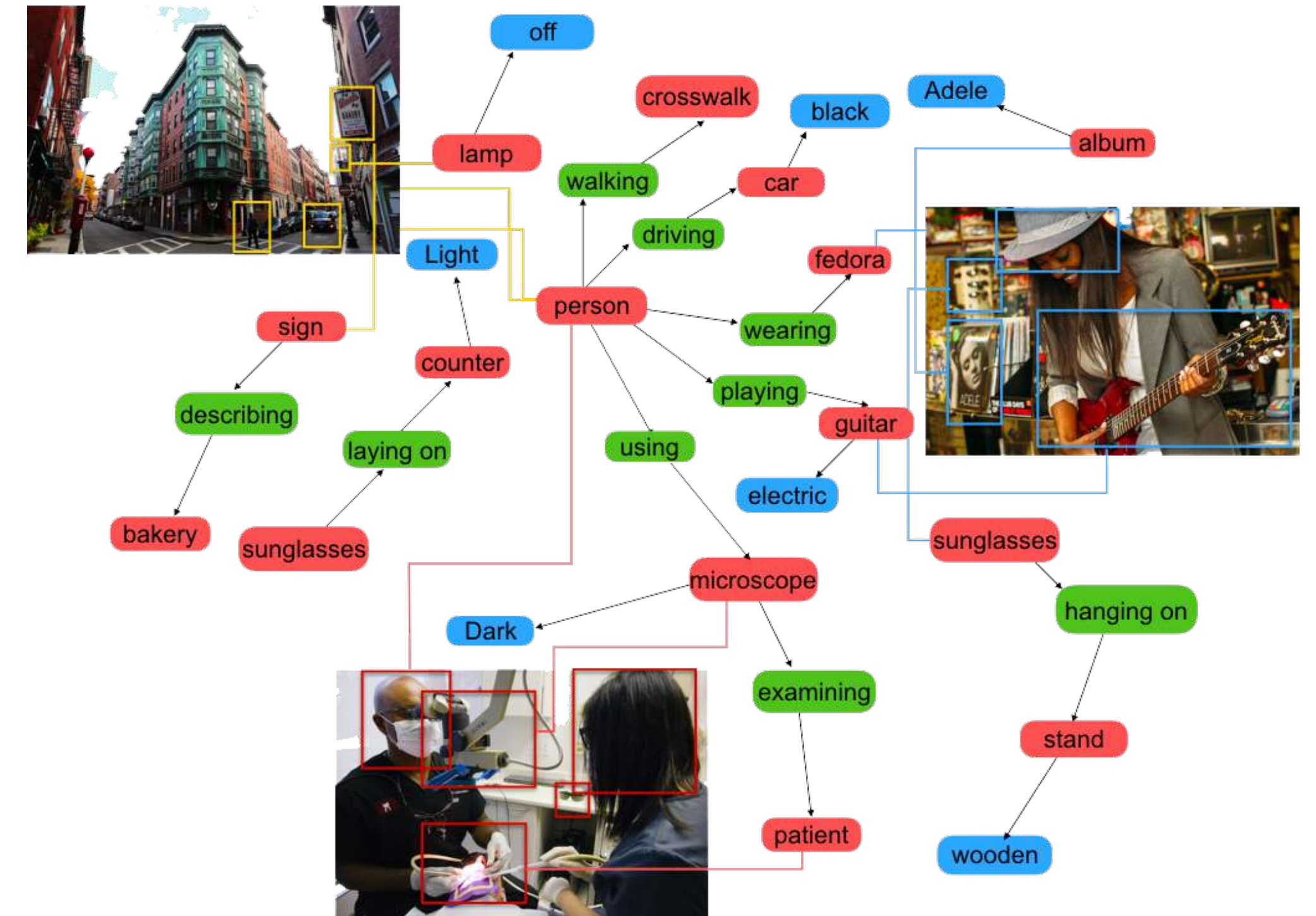Large-scale web videos & crowd teleoperation from online users

❖ How can robots take advantage of such knowledge?

Machine learning algorithms, e.g., meta-learning & imitation learning

## Open Question:

How to integrate procedural knowledge and declarative knowledge into a unified knowledge ontology for building intelligent algorithms in robotics?

# Acknowledgements

Declarative Knowledge ("That-Is")

Understanding the World

Robotics

Interacting with the World

Procedural Knowledge ("How-To")

http://ai.stanford.edu/~yukez/

yukez@cs.stanford.edu

Open Question:

How to integrate procedural knowledge and

declarative knowledge into a unified knowledge

ontology for building intelligent algorithms in

robotics?