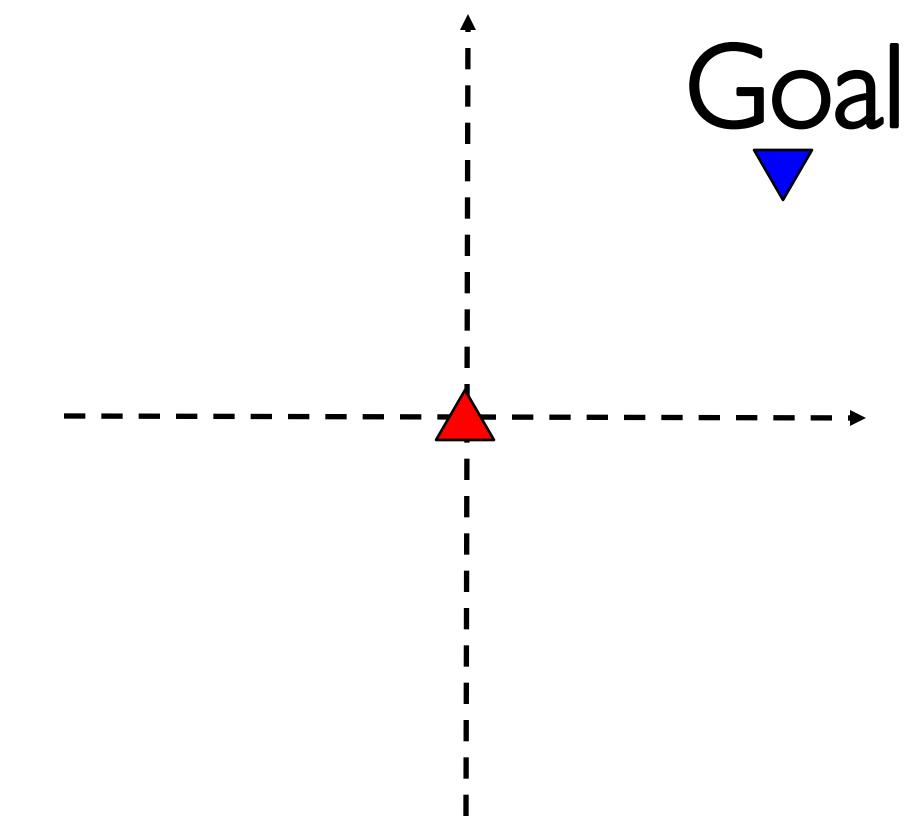




Robot with a first person camera



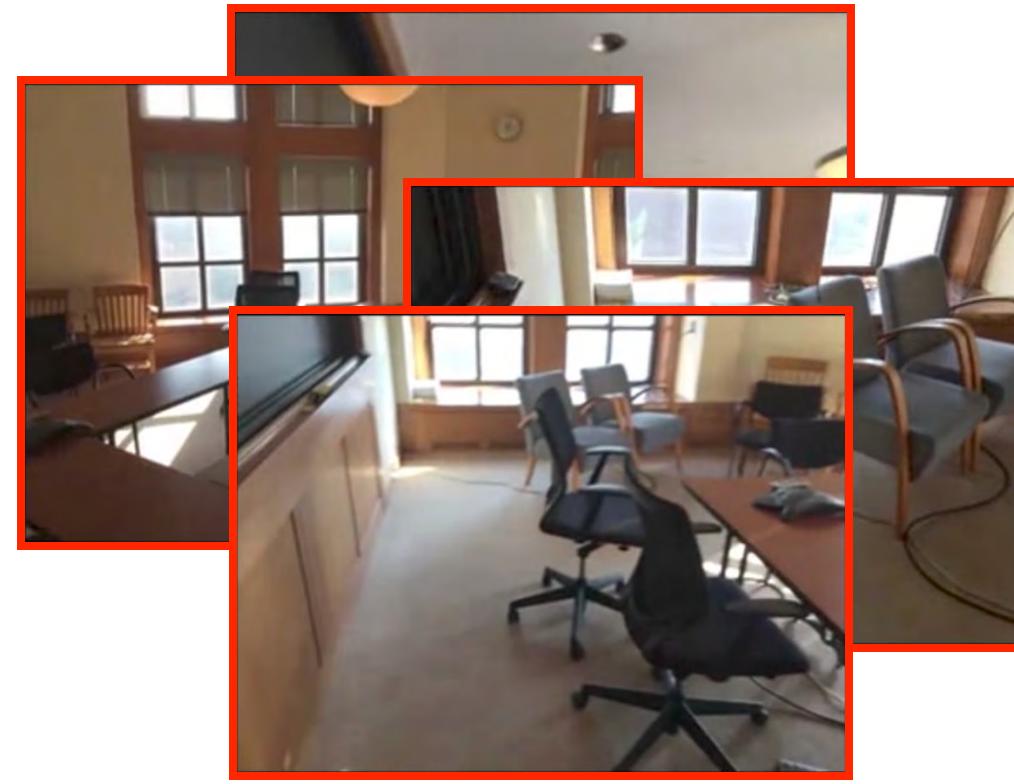
Dropped into a novel environment



**“Go
300 feet North,
400 feet East”**

**“Go Find a
Chair”**

Navigate around



Observed Images

Mapping

Planning

Hartley and Zisserman. 2000. Multiple View
Geometry in Computer Vision
Thrun, Burgard, Fox. 2005. Probabilistic Robotics

Canny. 1988. The complexity of robot motion
planning.

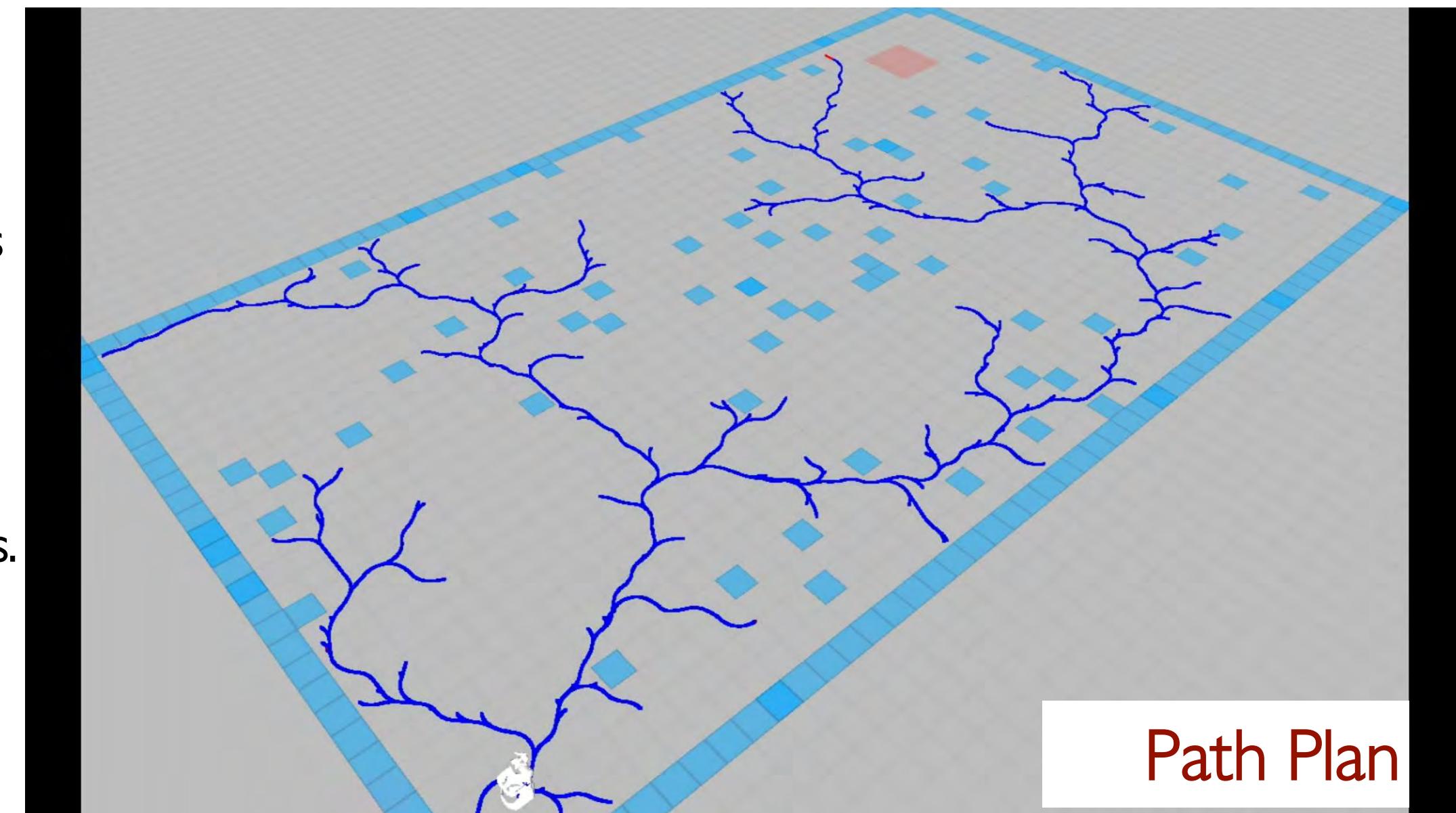
Kavraki et al. RA1996. Probabilistic roadmaps for
path planning in high-dimensional configuration spaces.

Lavalle and Kuffner. 2000. Rapidly-exploring
random trees: Progress and prospects.

Video Credits: Mur-Artal et al., Palmieri et al.

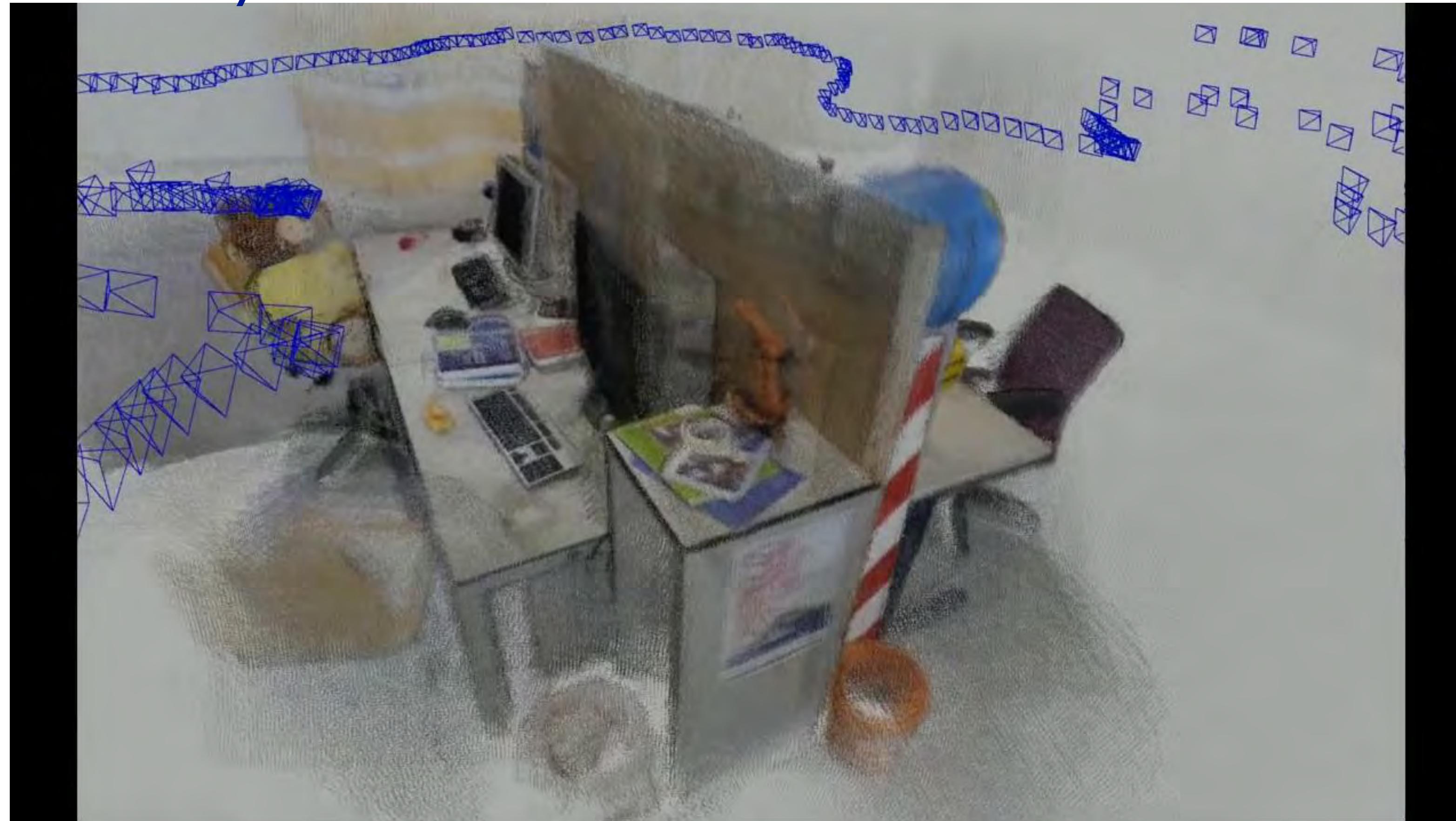


Geometric Reconstruction



Path Plan

Unnecessary

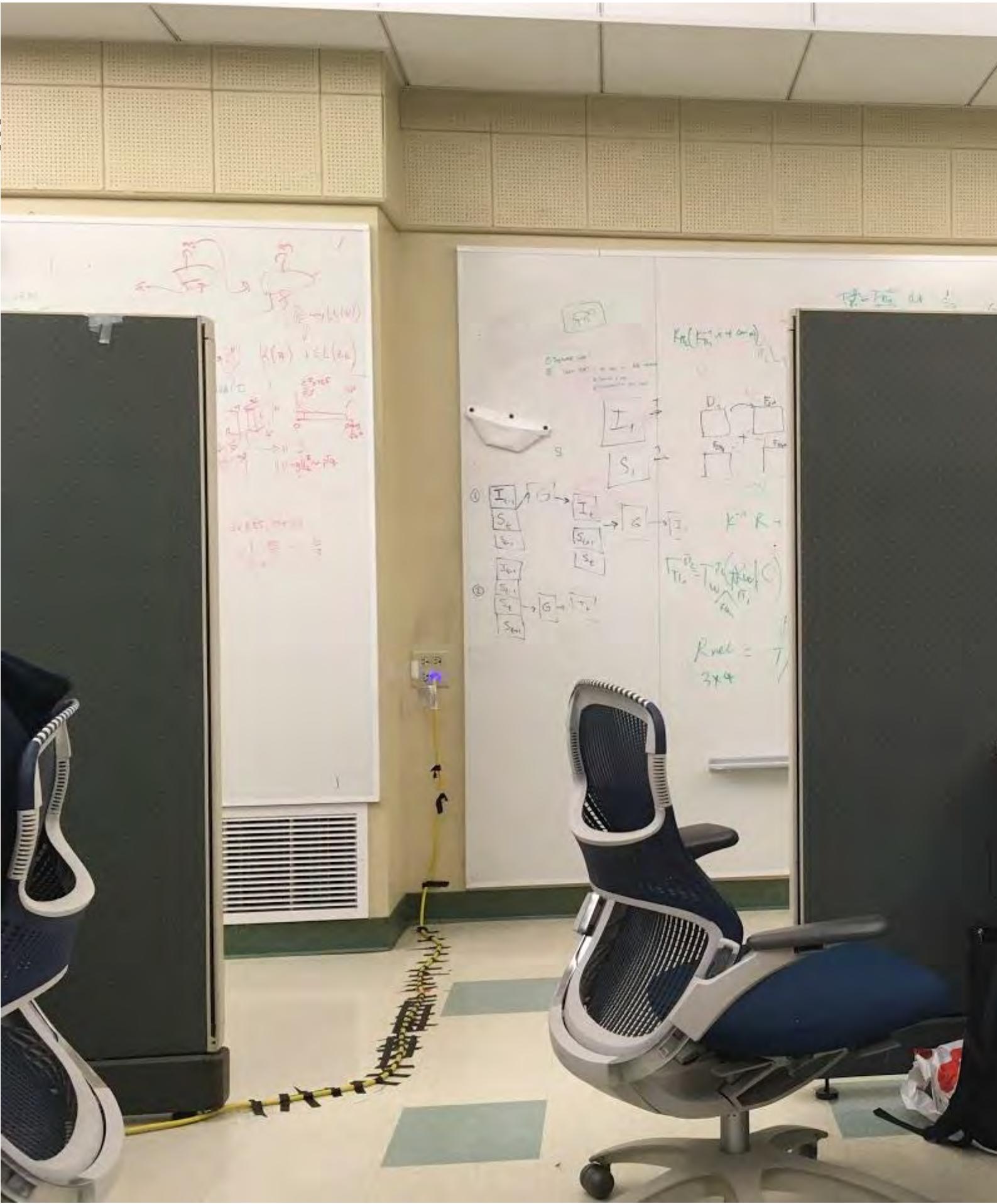


Do we need to tediously reconstruct everything on this table?

Video Credit: Mur-Artal and Tardos, TRobotics 2016. ORB-SLAM2: an Open-Source SLAM System for Monocular, Stereo and RGB-D Cameras.

Insufficient

Geometric World



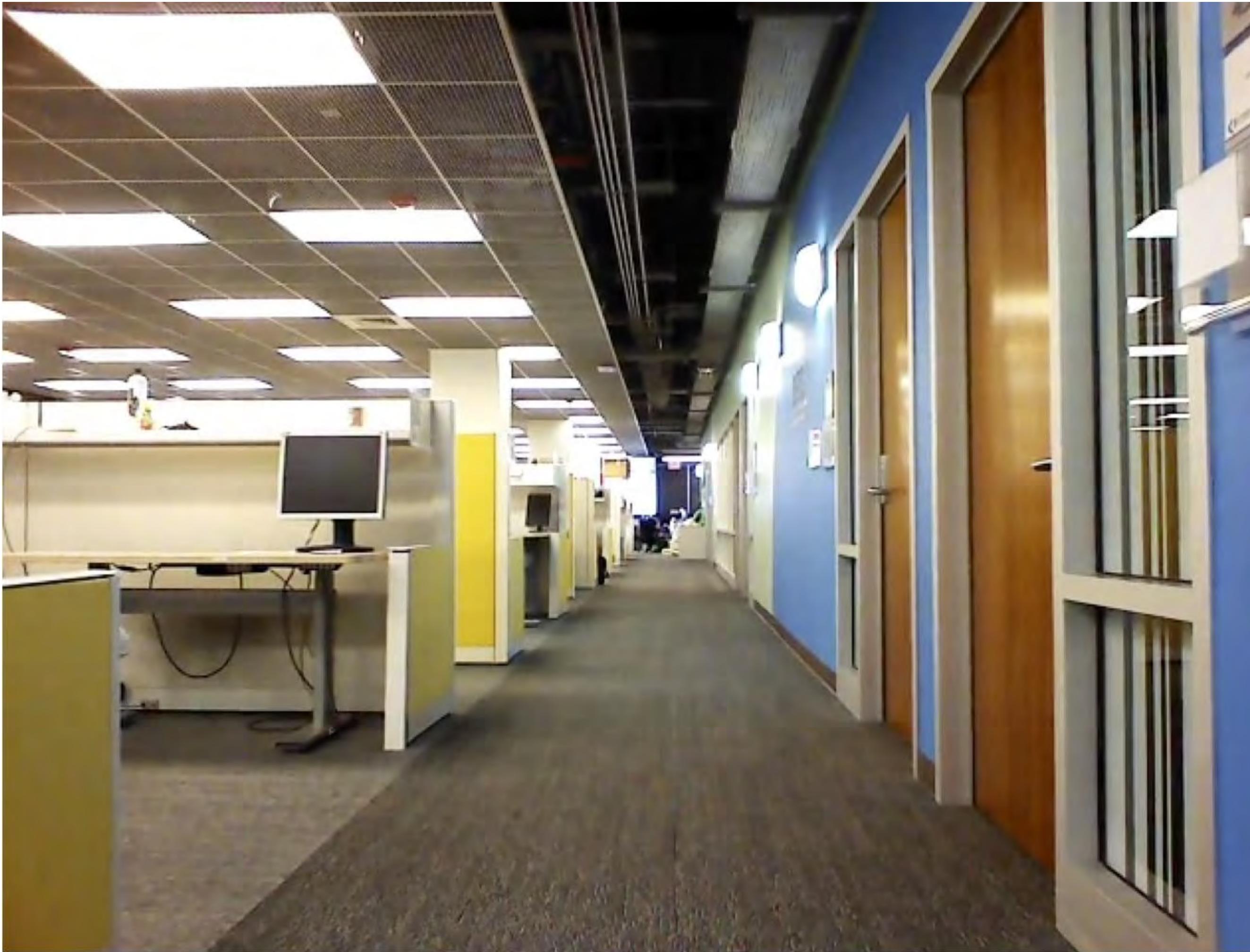
Can't speculate about space not directly observed.

Insufficient



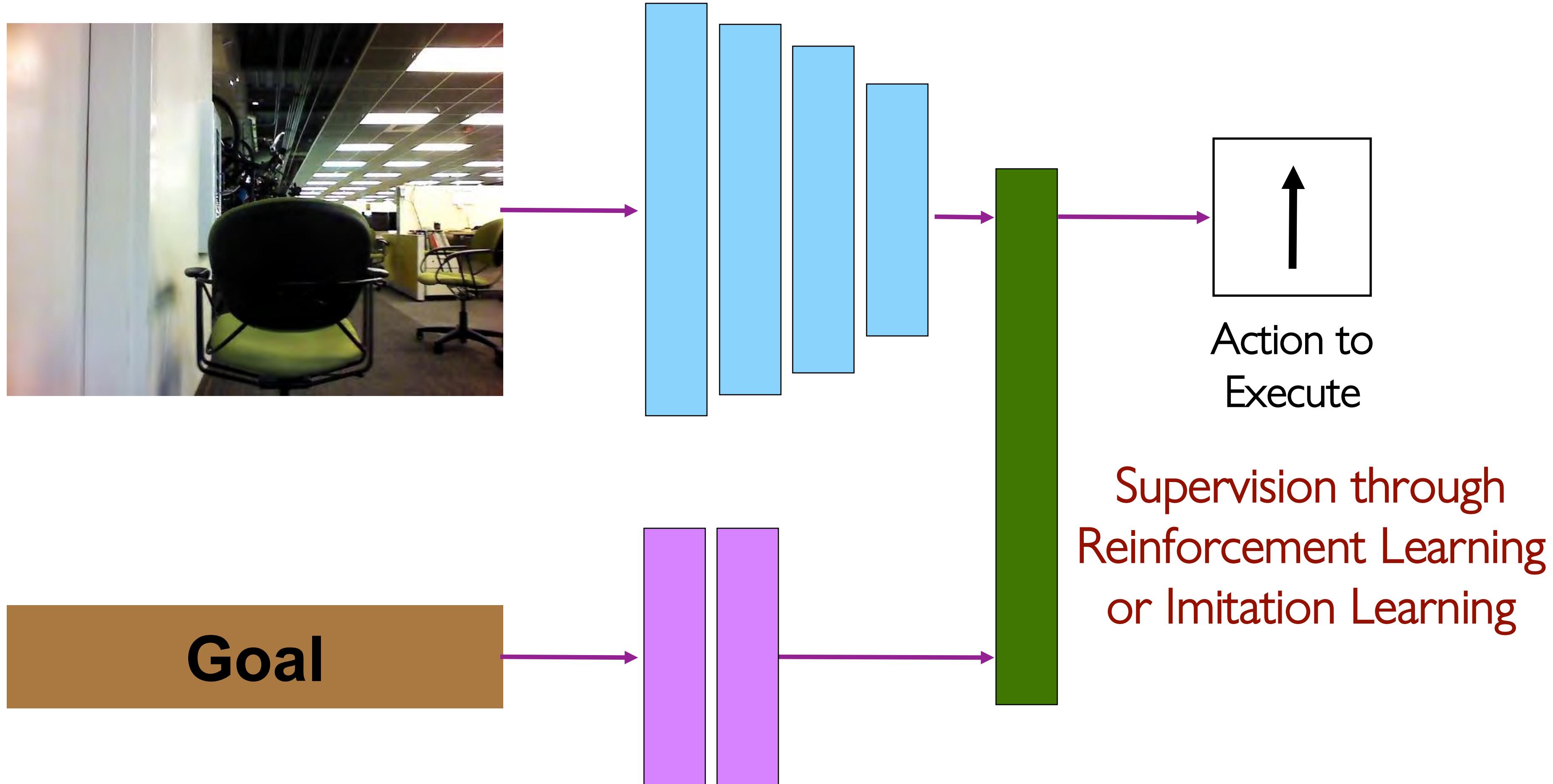
Can't exploit patterns in layout of indoor spaces.

Insufficient



Ignore navigational affordances.

Modern End-to-End Learned Navigation



Mnih et al., Nature 2014. Human-level control through deep reinforcement learning.

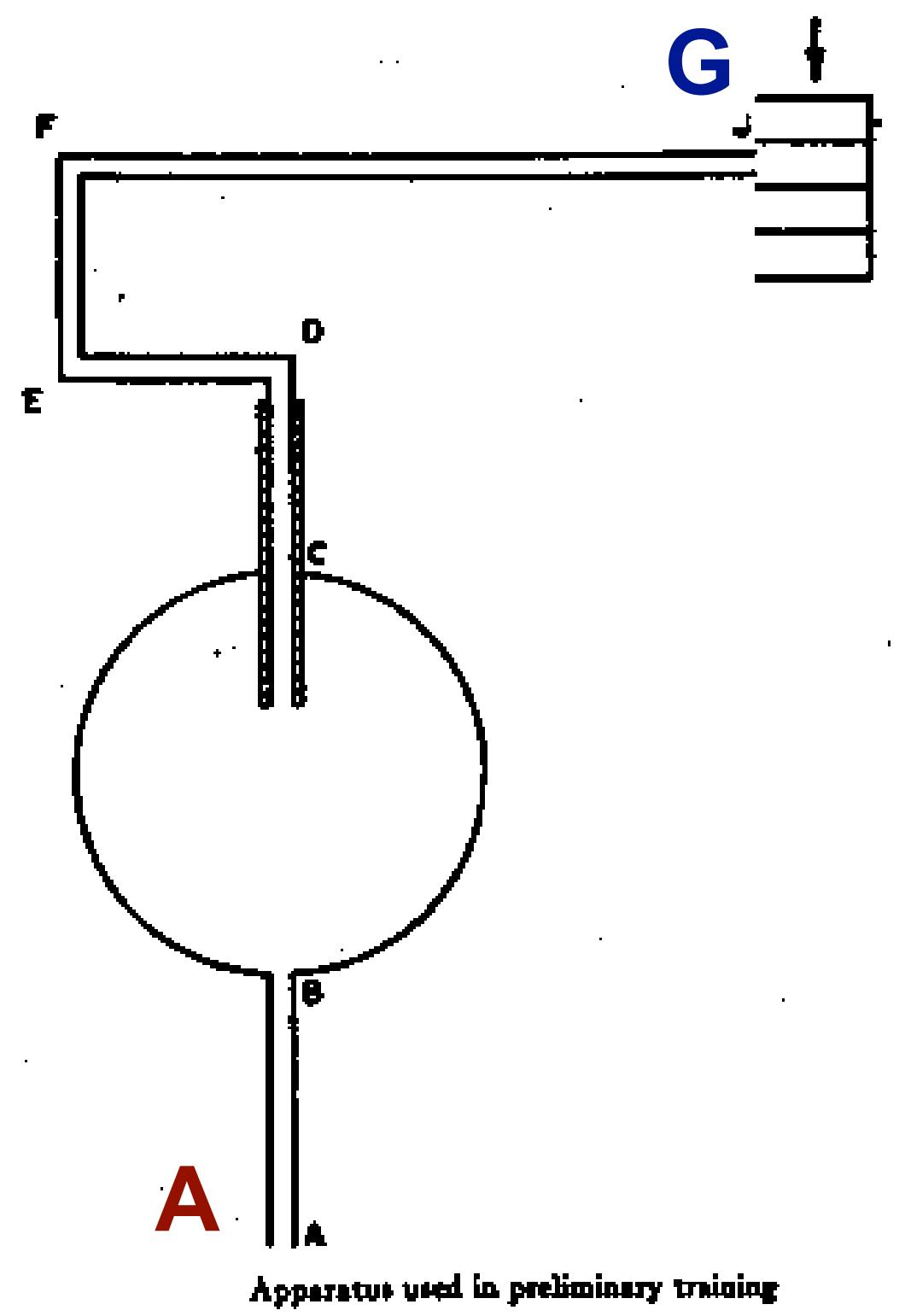
Levine et al., JMLR 2015. End-to-End Training of Deep Visuomotor Policies.

Zhu et al., ICRA 2017. Target-driven Visual Navigation in Indoor Scenes using Deep Reinforcement Learning.

COGNITIVE MAPS IN RATS AND MEN¹

BY EDWARD C. TOLMAN

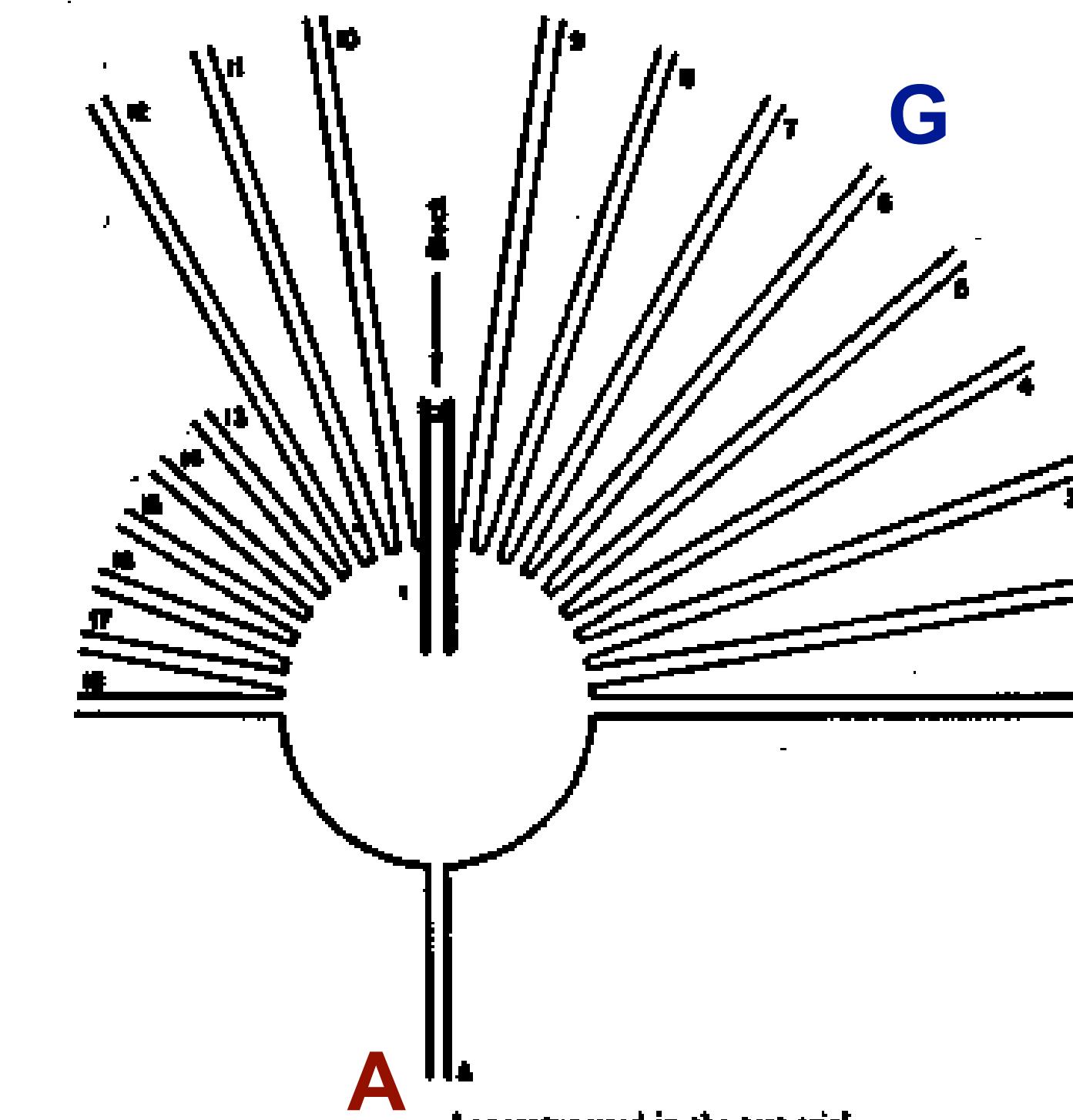
University of California



Apparatus used in preliminary training

FIG. 15

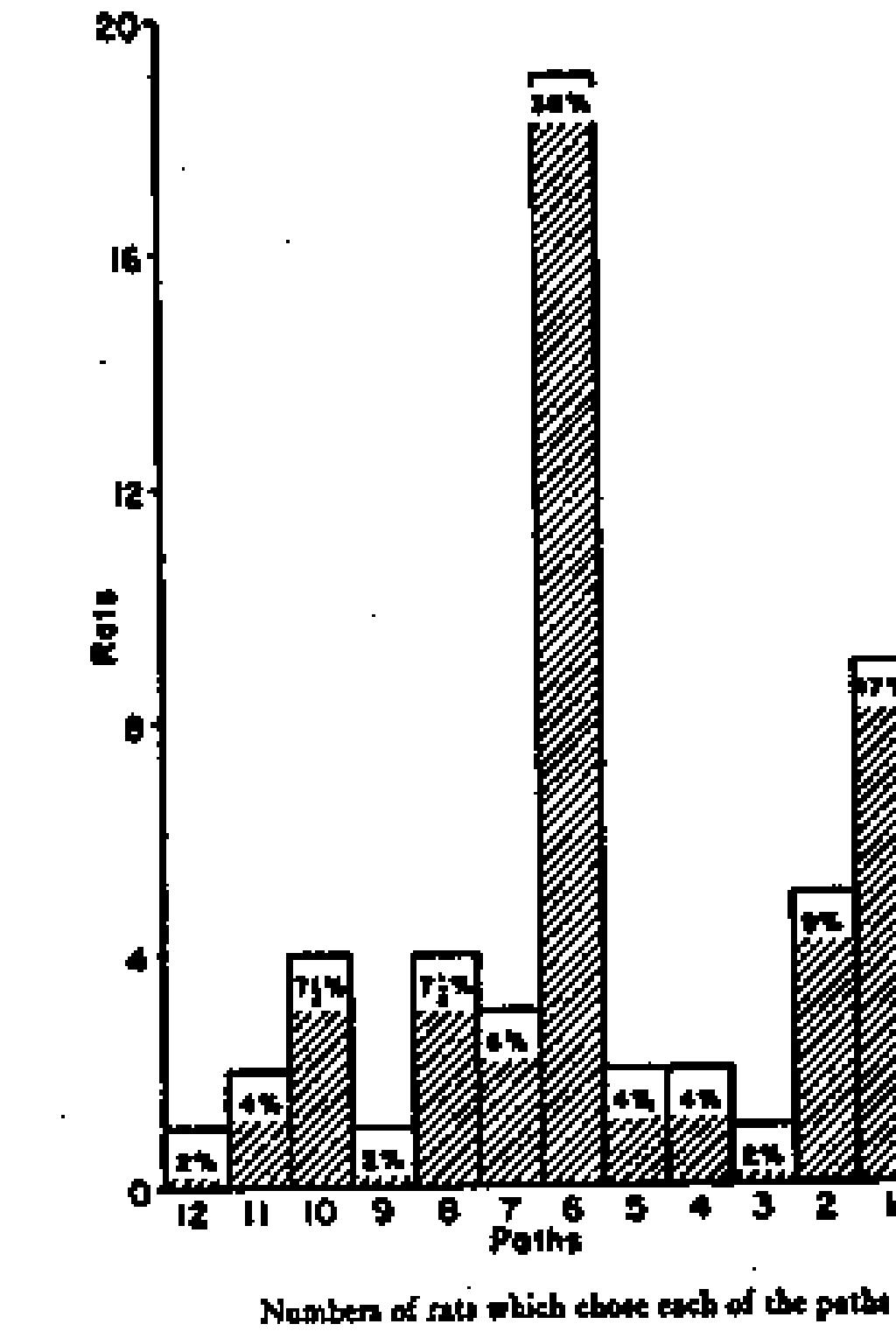
(From E. C. Tolman, B. F. Ritchie and D. Kalish, Studies in spatial learning. I. Orientation and the short-cut. *J. exp. Psychol.*, 1946, 36, p. 16.)



Apparatus used in the test trial

FIG. 16

(From E. C. Tolman, B. F. Ritchie and D. Kalish, Studies in spatial learning. I. Orientation and short-cut. *J. exp. Psychol.*, 1946, 36, p. 17.)



Numbers of rats which chose each of the paths

FIG. 17

(From E. C. Tolman, B. F. Ritchie and D. Kalish, Studies in spatial learning. I. Orientation and the short-cut. *J. exp. Psychol.*, 1946, 36, p. 19.)

Navigating to Objects in the Real World



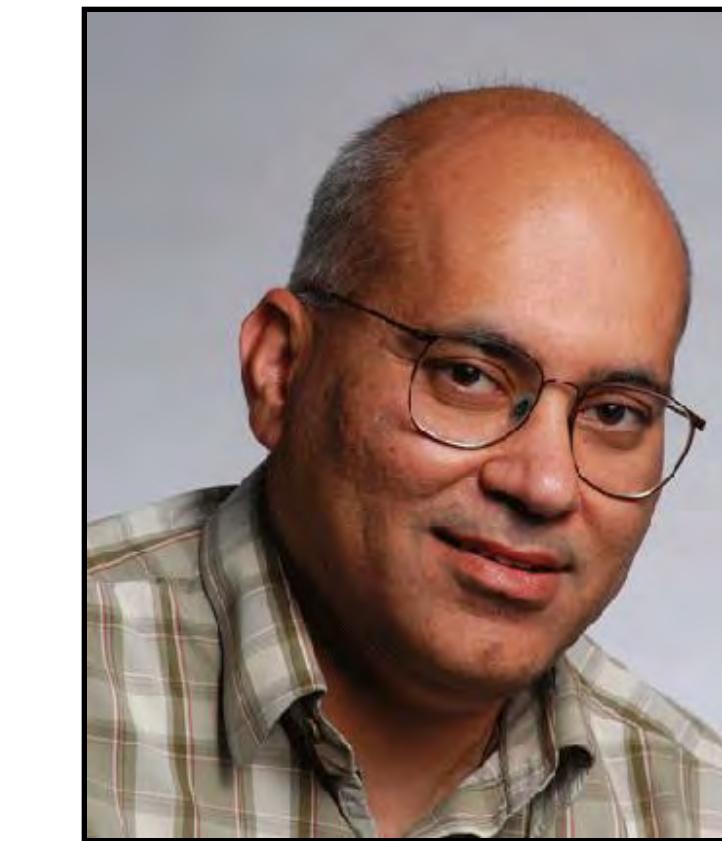
Theophile Gervet ¹



Soumith Chintala ⁴



Dhruv Batra ^{3,4}



Jitendra Malik ^{2,4}



Devendra Chaplot ⁴

¹ Carnegie
Mellon
University

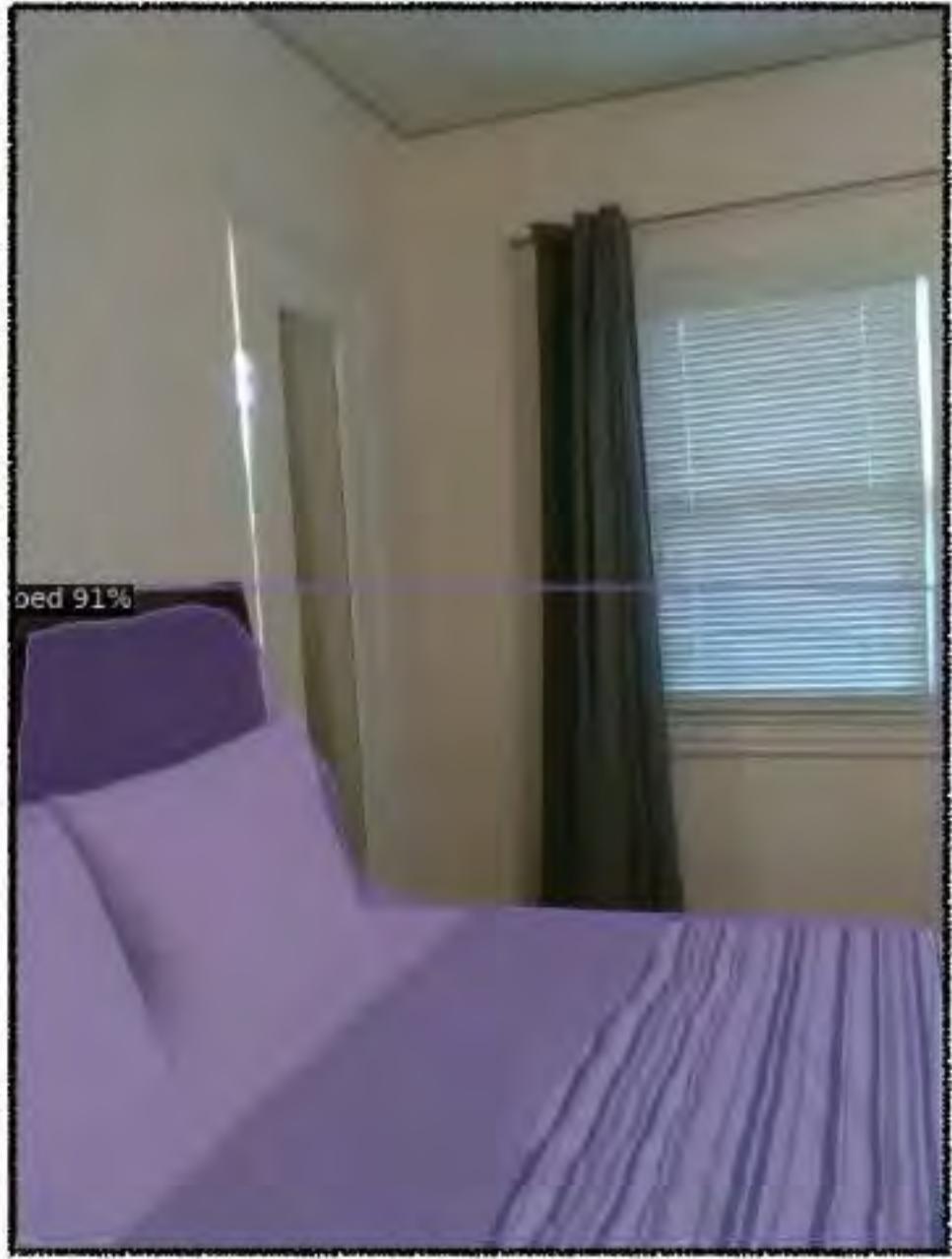


Berkeley
UNIVERSITY OF CALIFORNIA

³ Georgia
Tech

⁴ Meta AI

Observation



Goal: *plant*

Predicted
Semantic Map

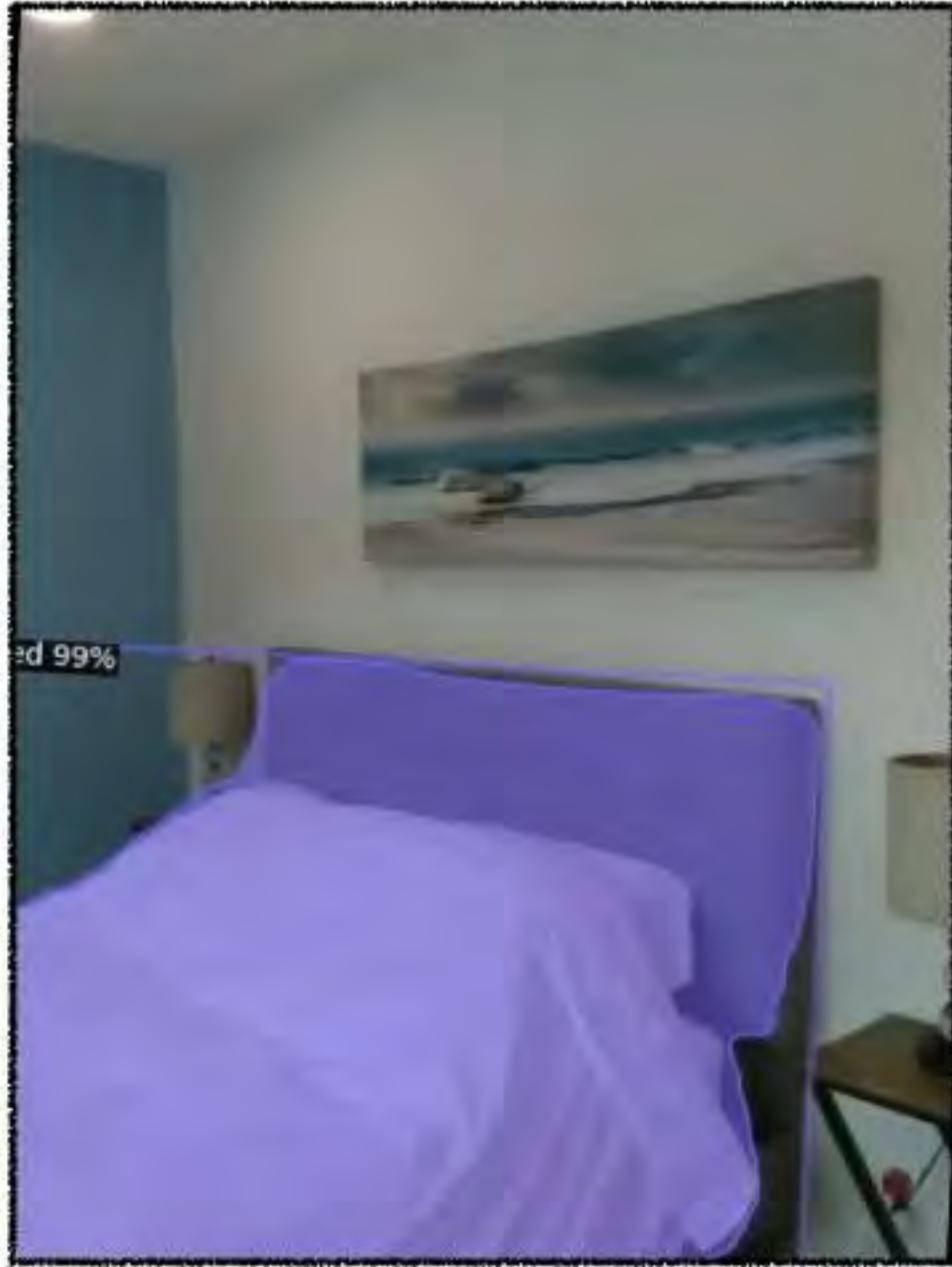


Third-person view



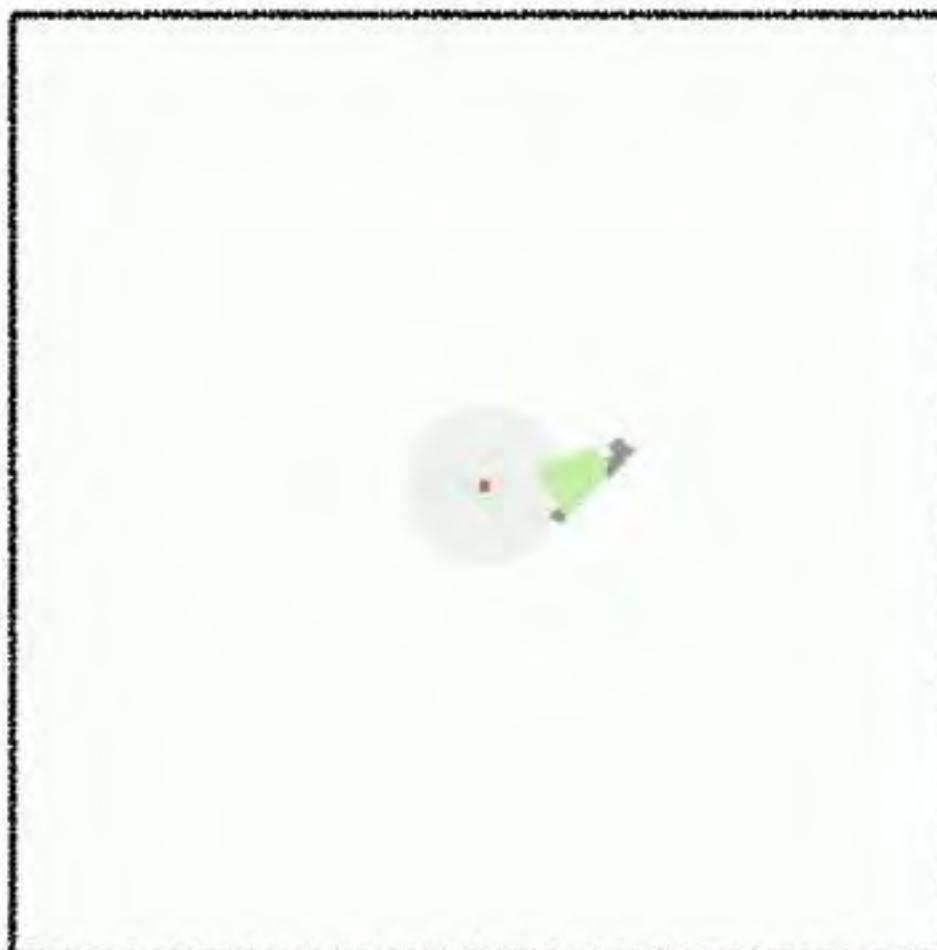
Navigable Area	3: bed	7: oven	11: clock
0: chair	4: toilet	8: sink	12: vase
1: couch	5: tv	9: refrigerator	13: cup
2: potted plant	6: dining-table	10: book	14: bottle

Observation



Goal: chair

Predicted
Semantic Map



Third-person view



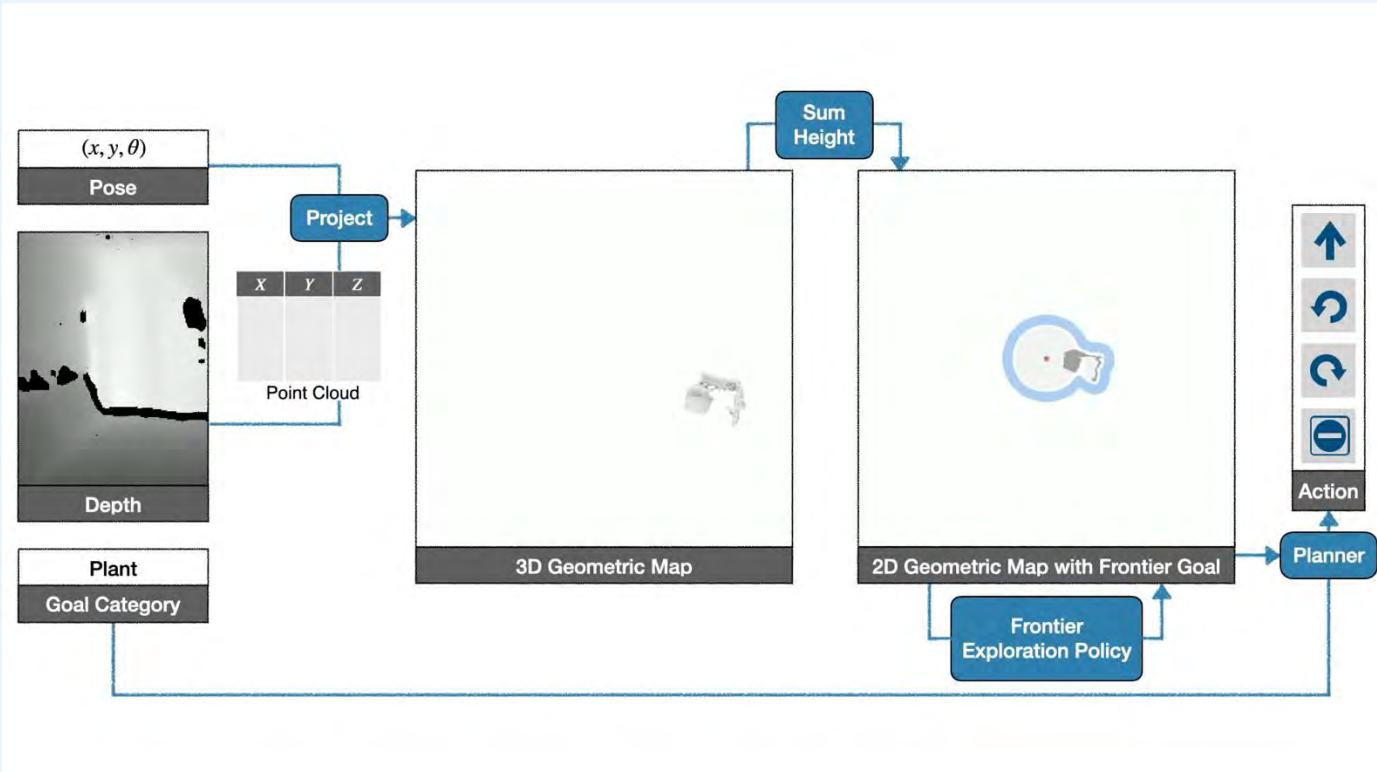
Navigable Area	3: bed	7: oven	11: clock
0: chair	4: toilet	8: sink	12: vase
1: couch	5: tv	9: refrigerator	13: cup
2: potted plant	6: dining-table	10: book	14: bottle



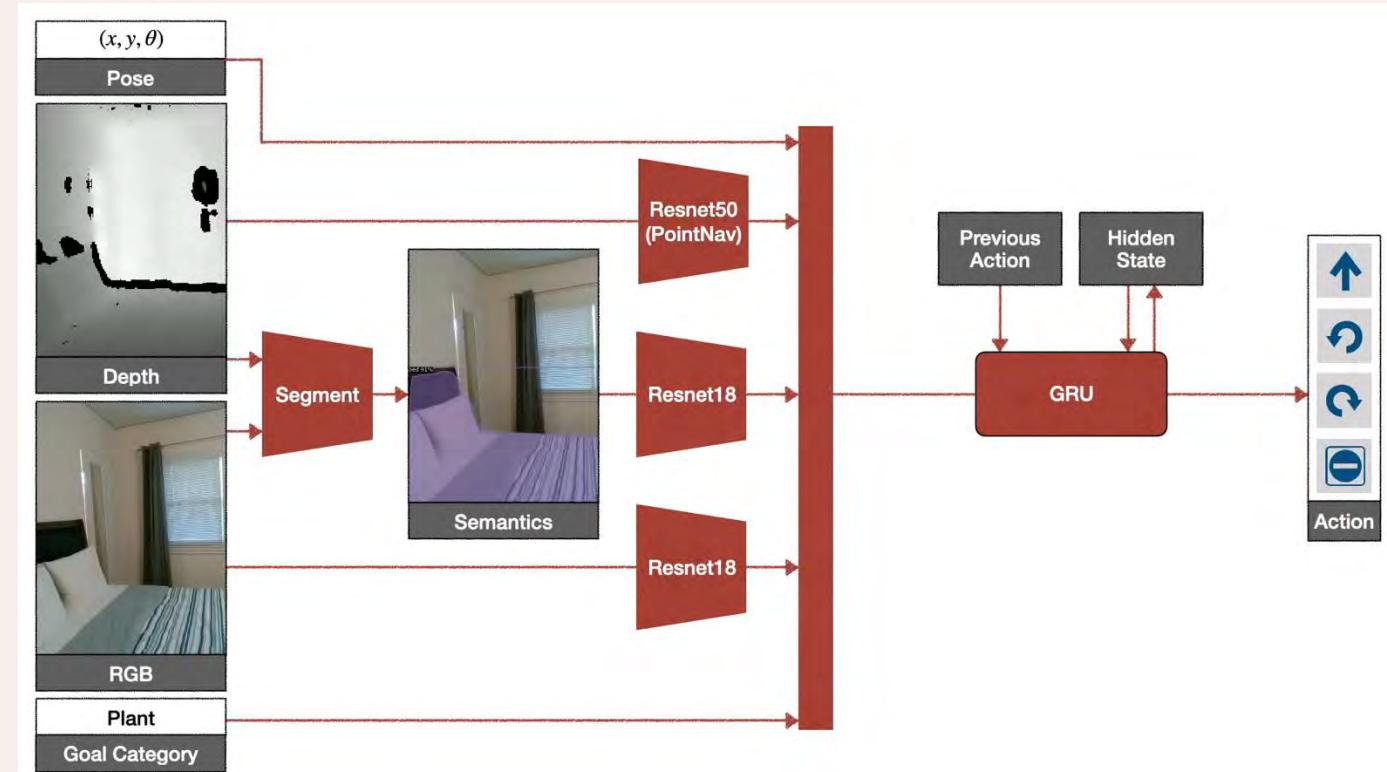
Empirical Evaluation
3 Approaches
6 Unseen Homes
6 Goal Object Categories

Methods

Classical

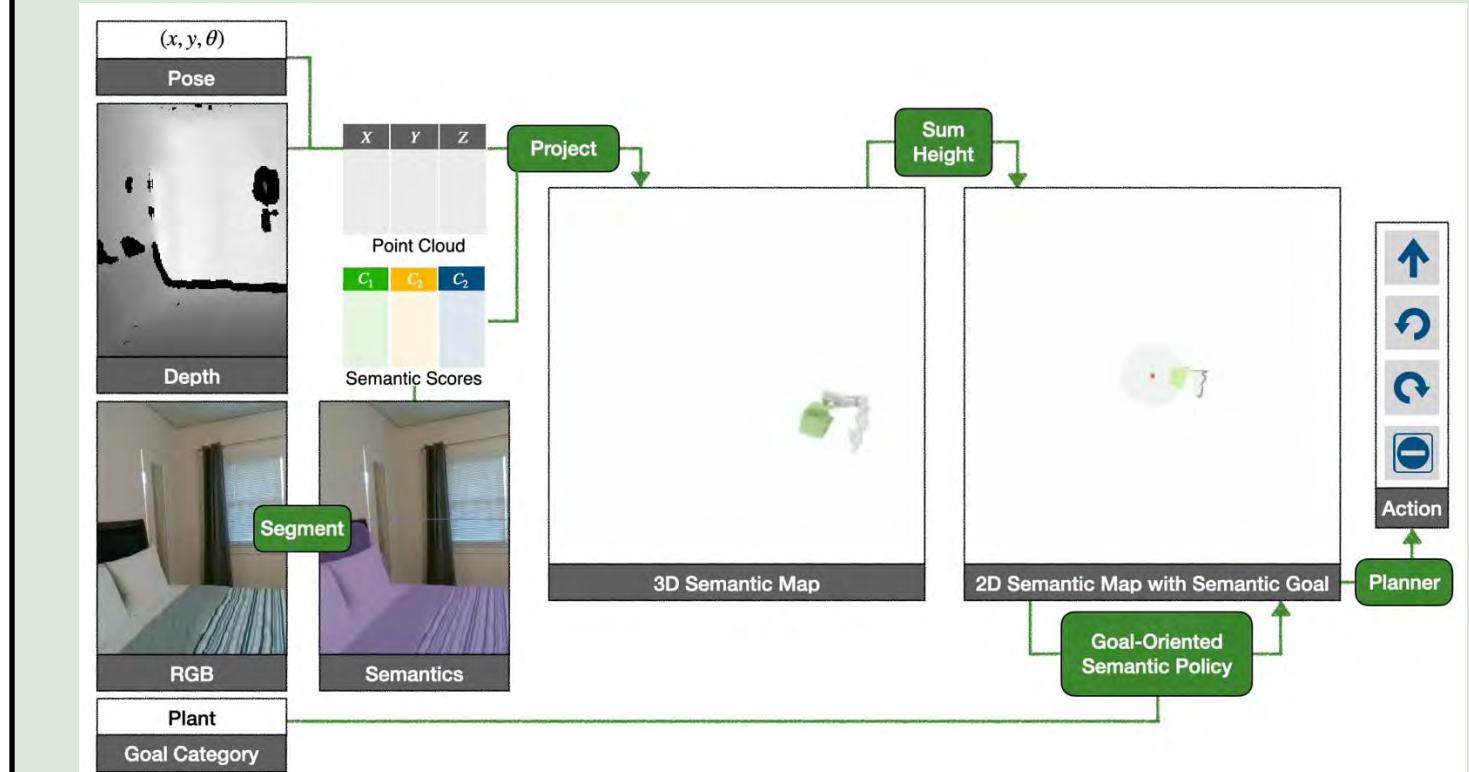


End-to-end Learning



- ▶ Geometric Map
- ▶ Heuristic Exploration
- ▶ No Training

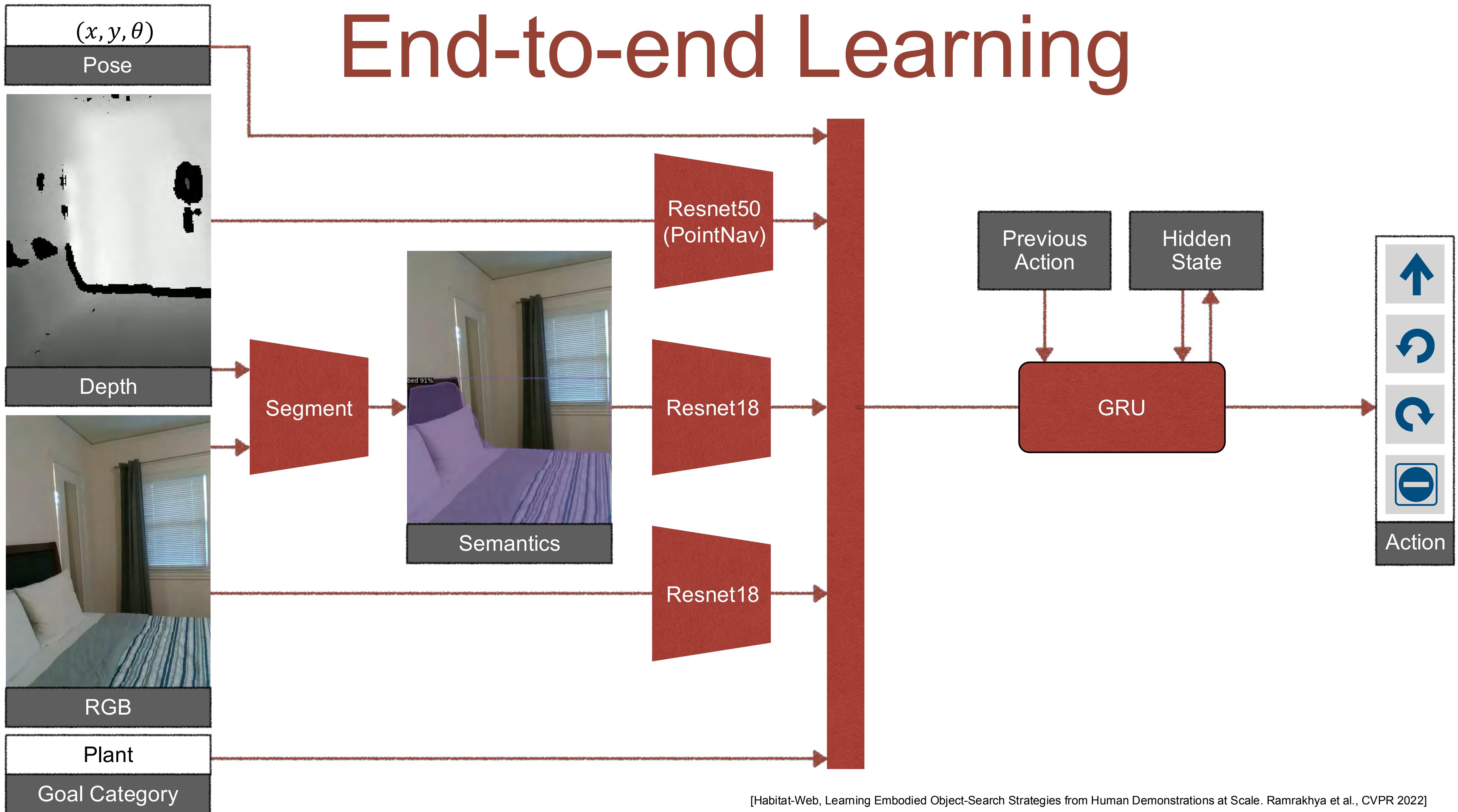
Modular Learning



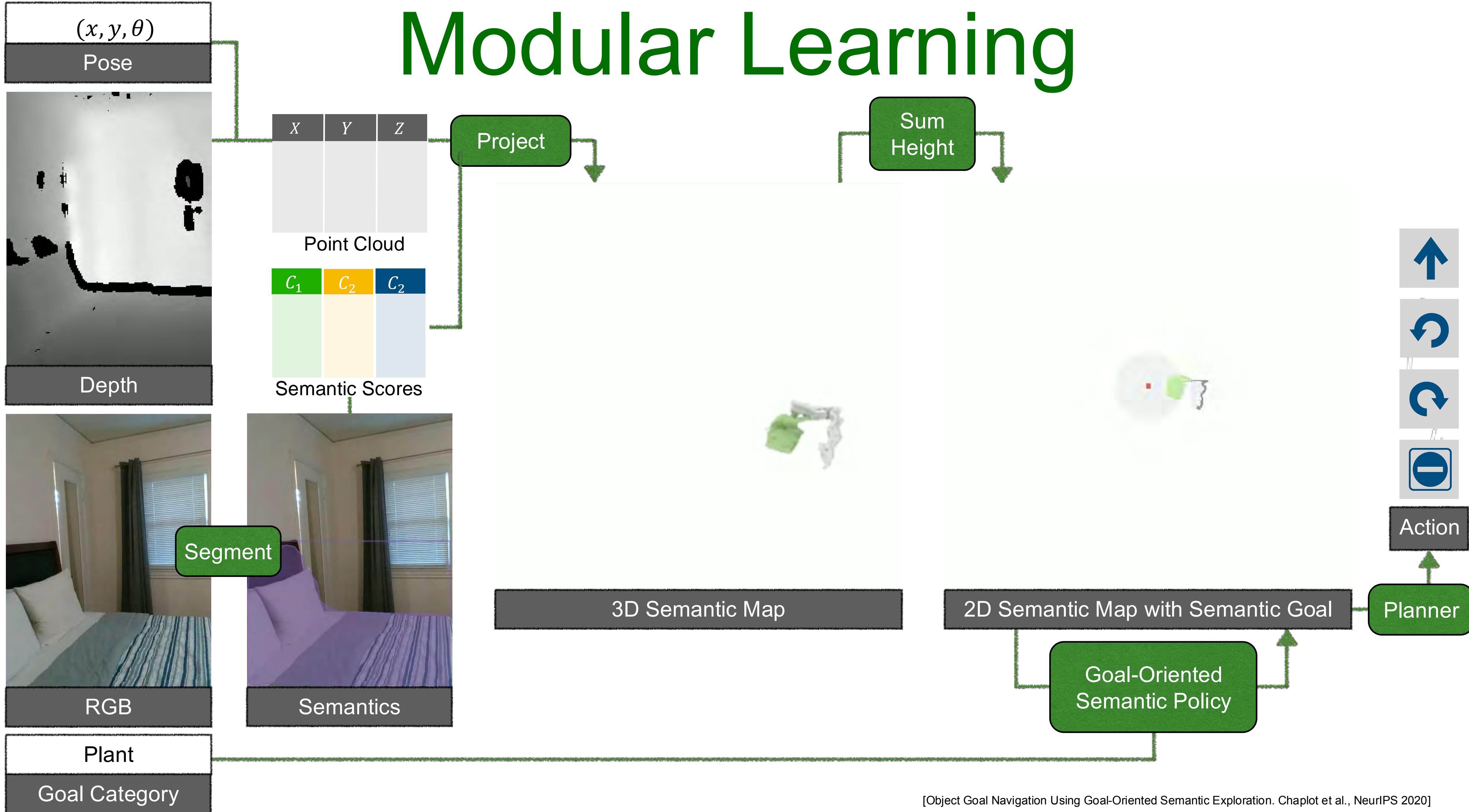
- ▶ End-to-end
- ▶ Large-scale IL + RL fine-tuning
 - ▶ 77,000 human trajectories
 - ▶ 200M frames of RL

- ▶ Semantic Map
- ▶ Goal-Oriented Exploration
- ▶ 10M frames of RL

End-to-end Learning



Modular Learning



Habitat



AI2-Thor



Results

Success Rate

SPL

■ Sim

■ Real World

Modular Learning

Classical

End-to-end Learning

0.00

0.25

0.50

0.75

1.00

Goal: couch

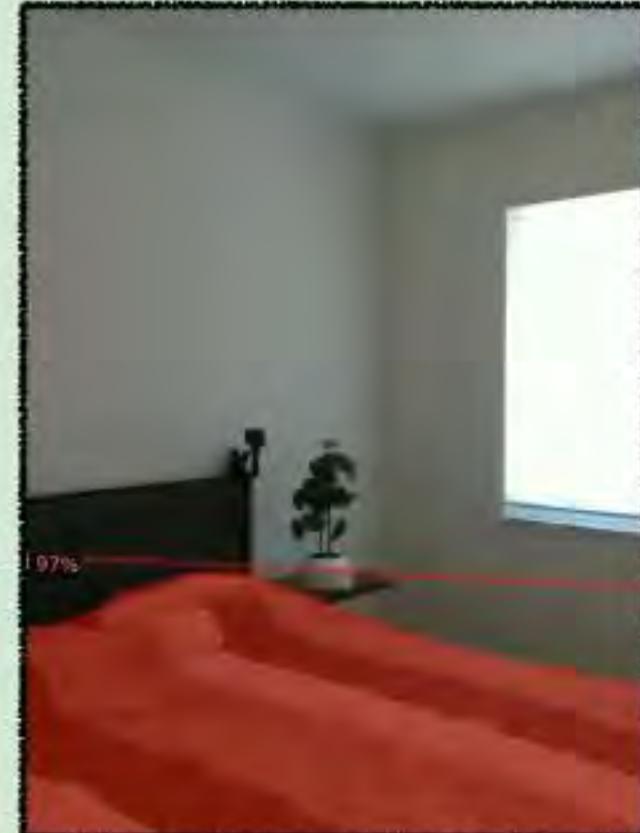
SPL: 0.74, 78 steps

Modular

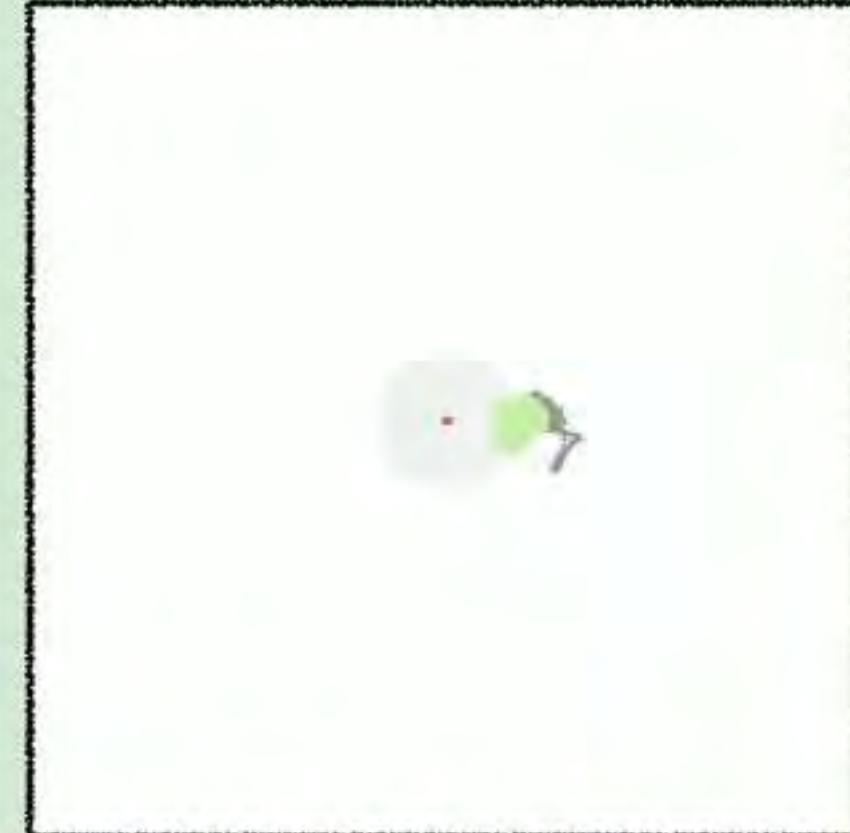
Third-person view



Success



Observation



Predicted
Semantic Map

SPL: 0.0, 121 steps

End-to-End

Third-person view



Failure



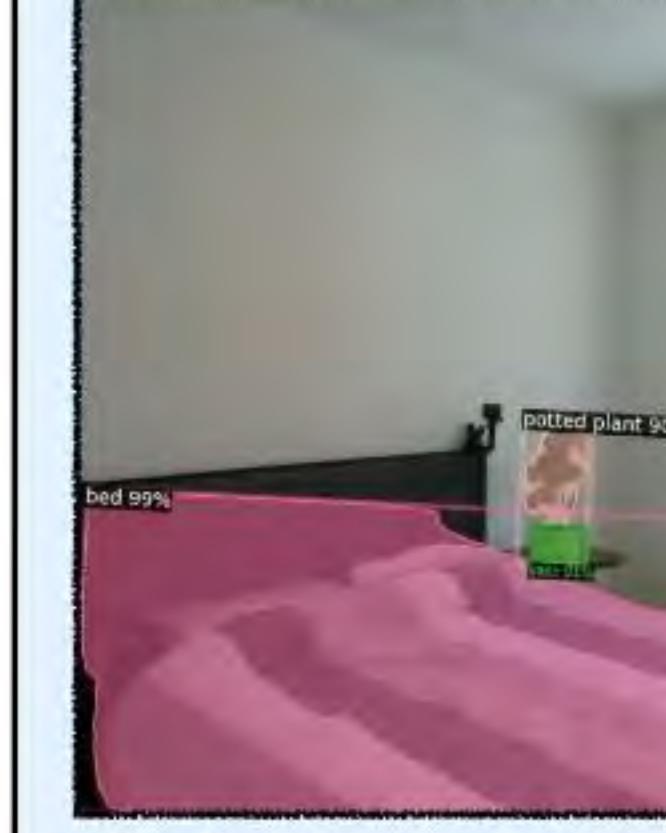
SPL: 0.33, 181 steps

Classical

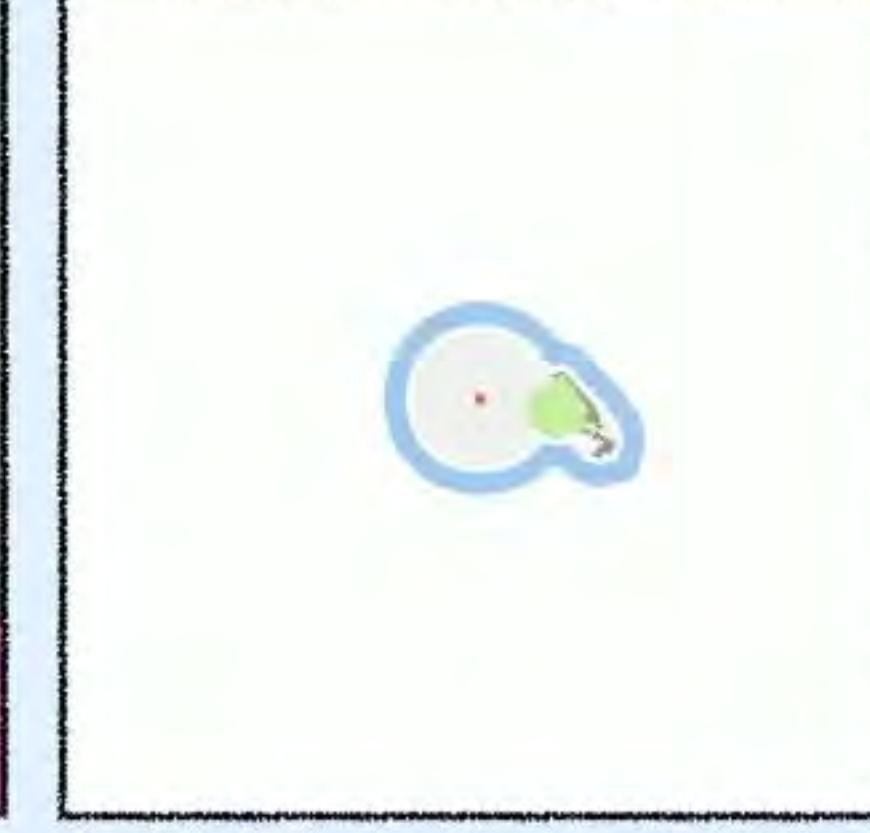
Third-person view



Success

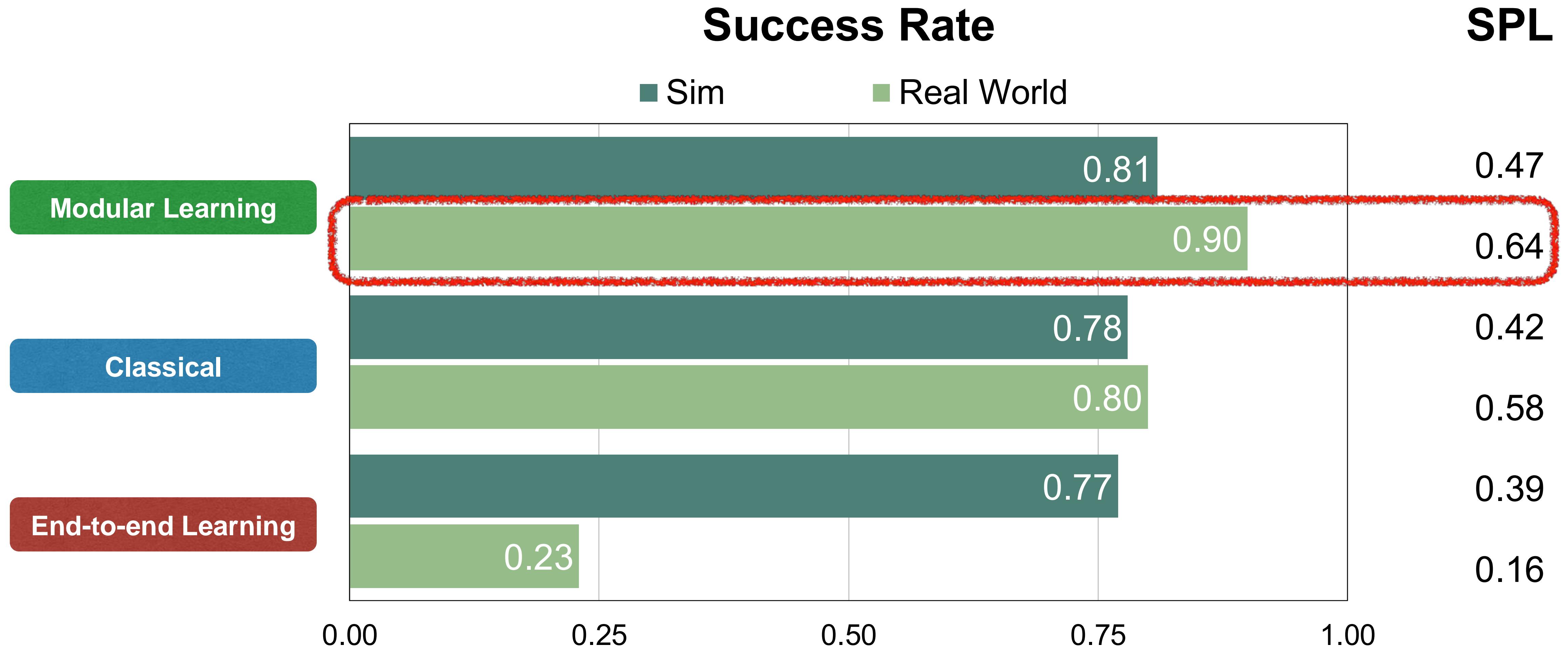


Observation



Predicted
Semantic Map

Modular Learning is Reliable



Observation



Goal: *toilet*

Predicted Semantic Map

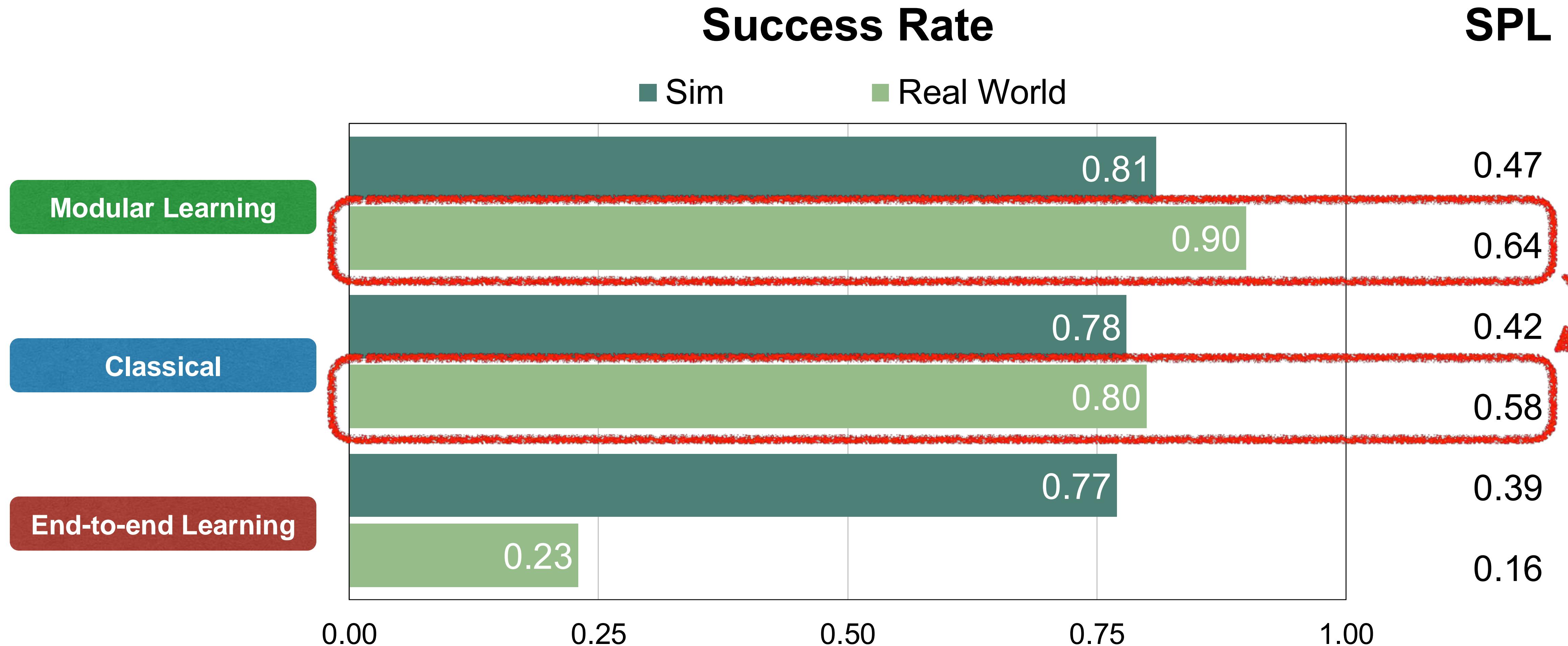


Third-person view



Navigable Area	3: bed	7: oven	11: clock
0: chair	4: toilet	8: sink	12: vase
1: couch	5: tv	9: refrigerator	13: cup
2: potted plant	6: dining-table	10: book	14: bottle

Classical vs Modular Learning



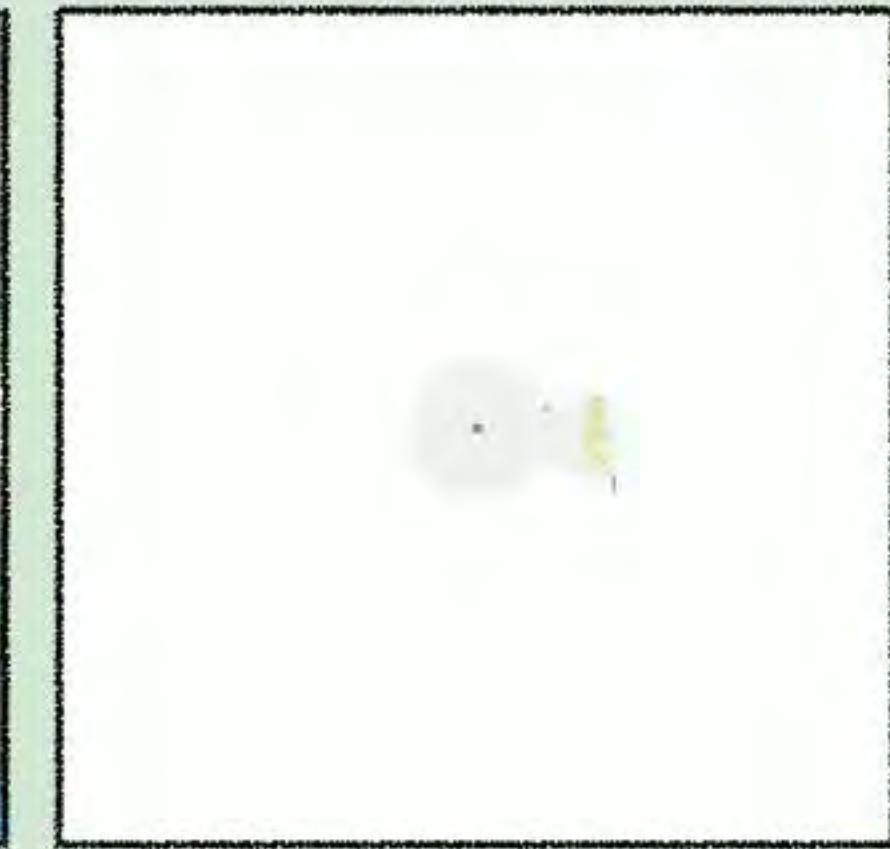
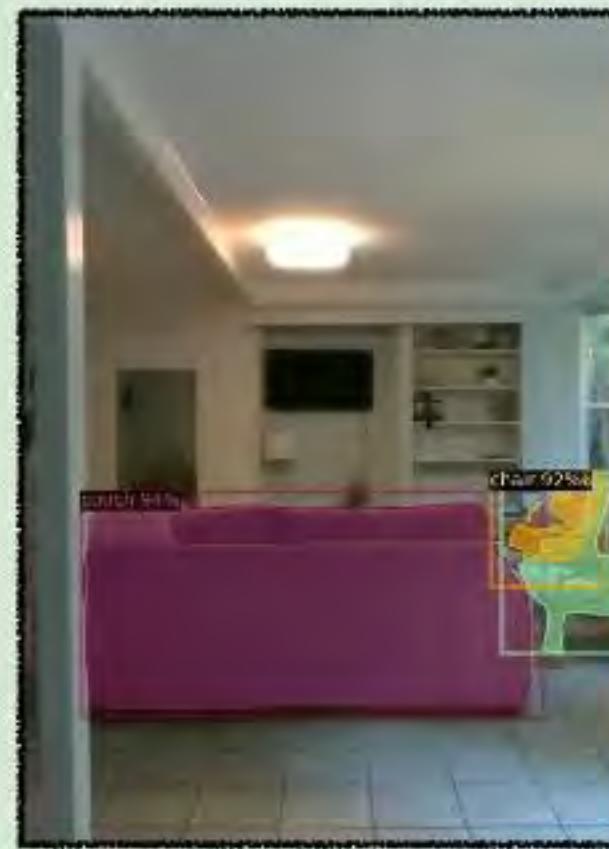
Classical vs Modular Learning

SPL: 0.90, 98 steps

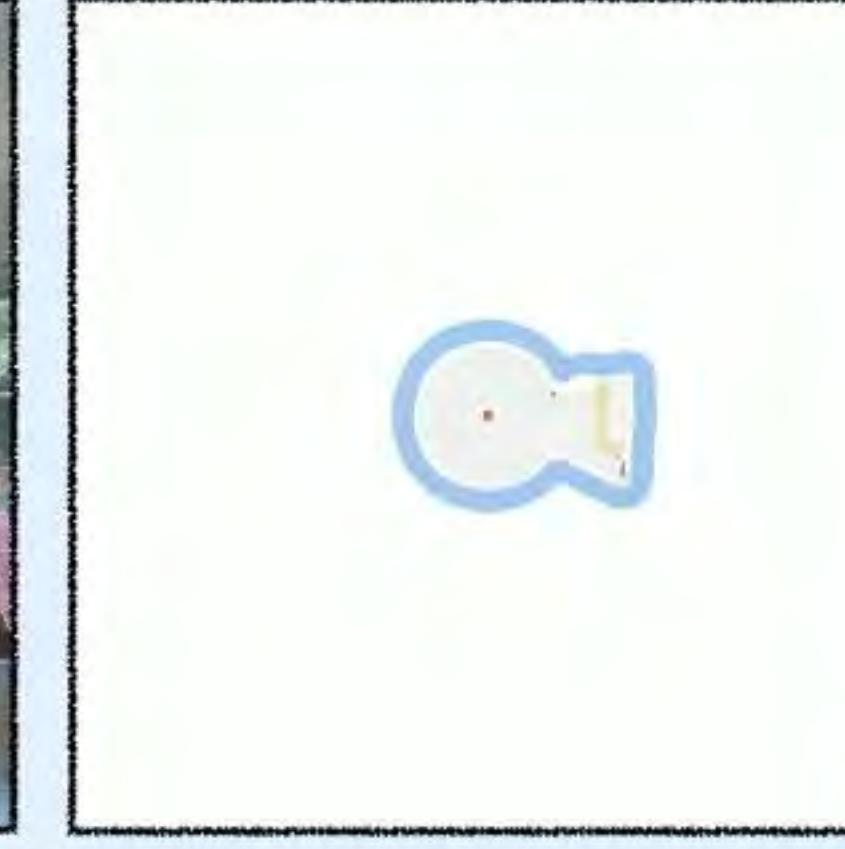
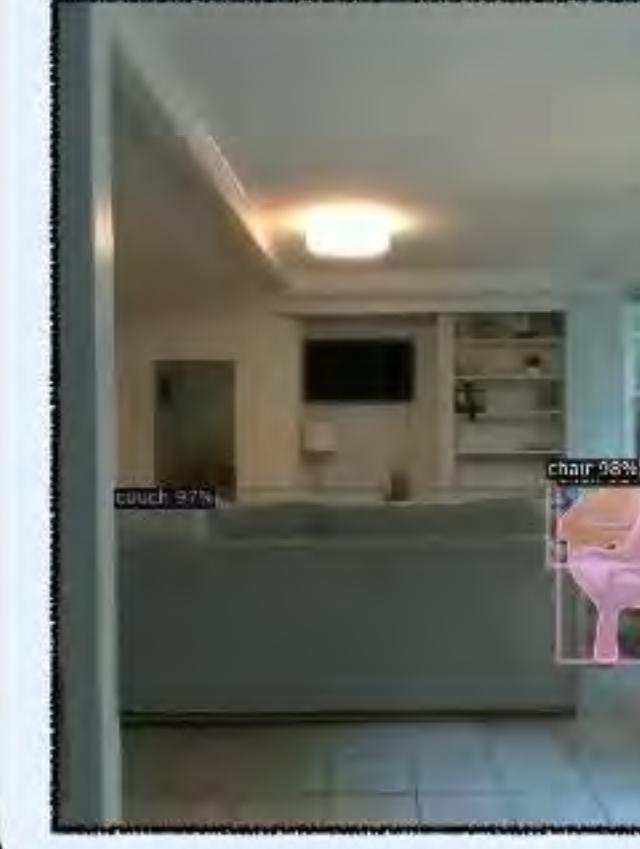
Goal: bed

SPL: 0.52, 152 steps

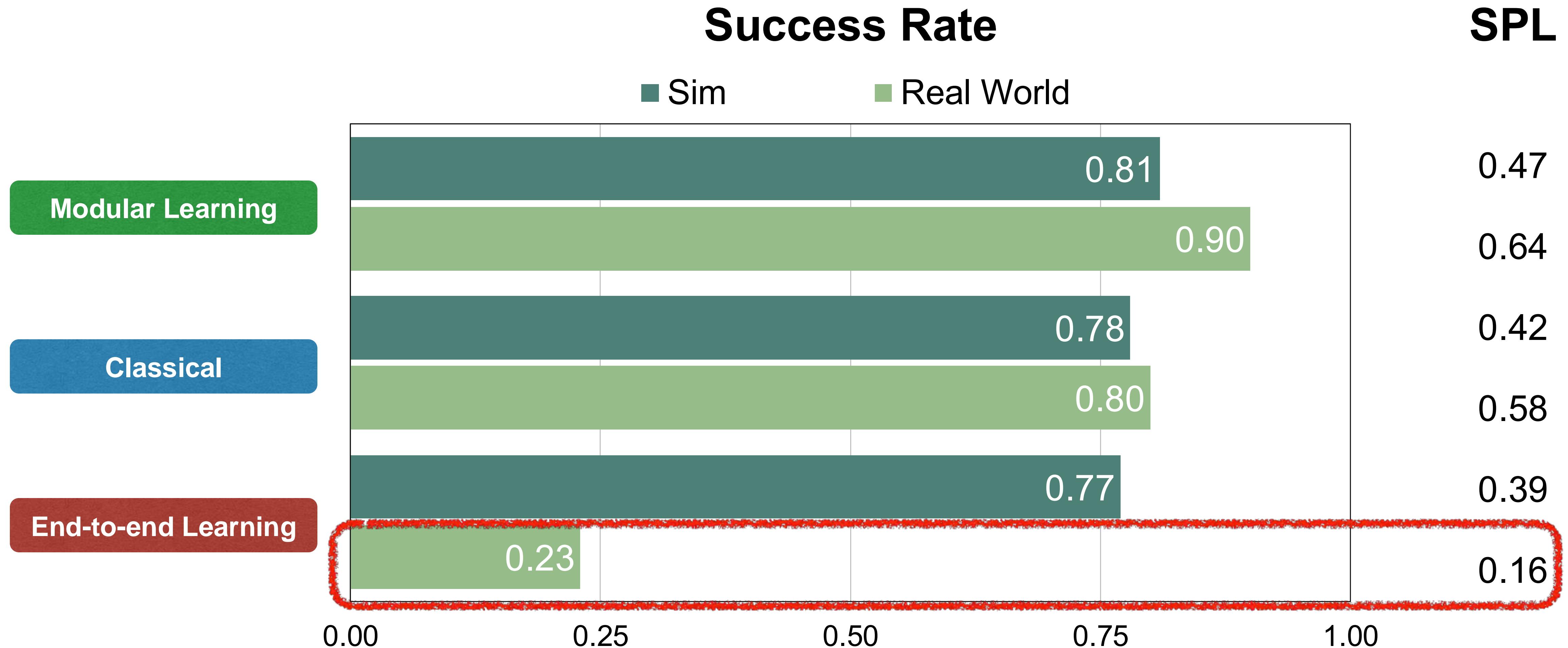
Semantic Exploration



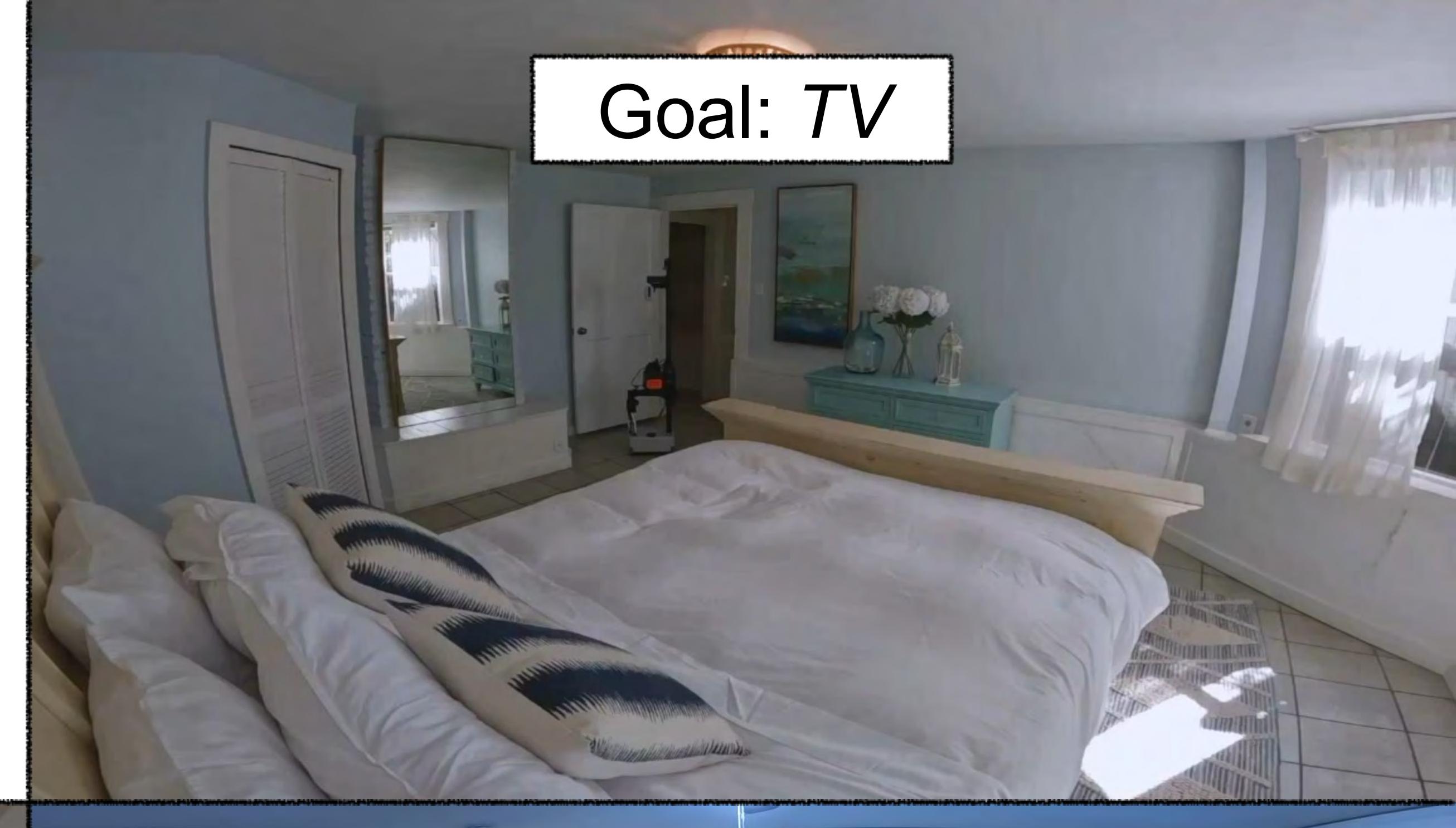
Frontier Exploration



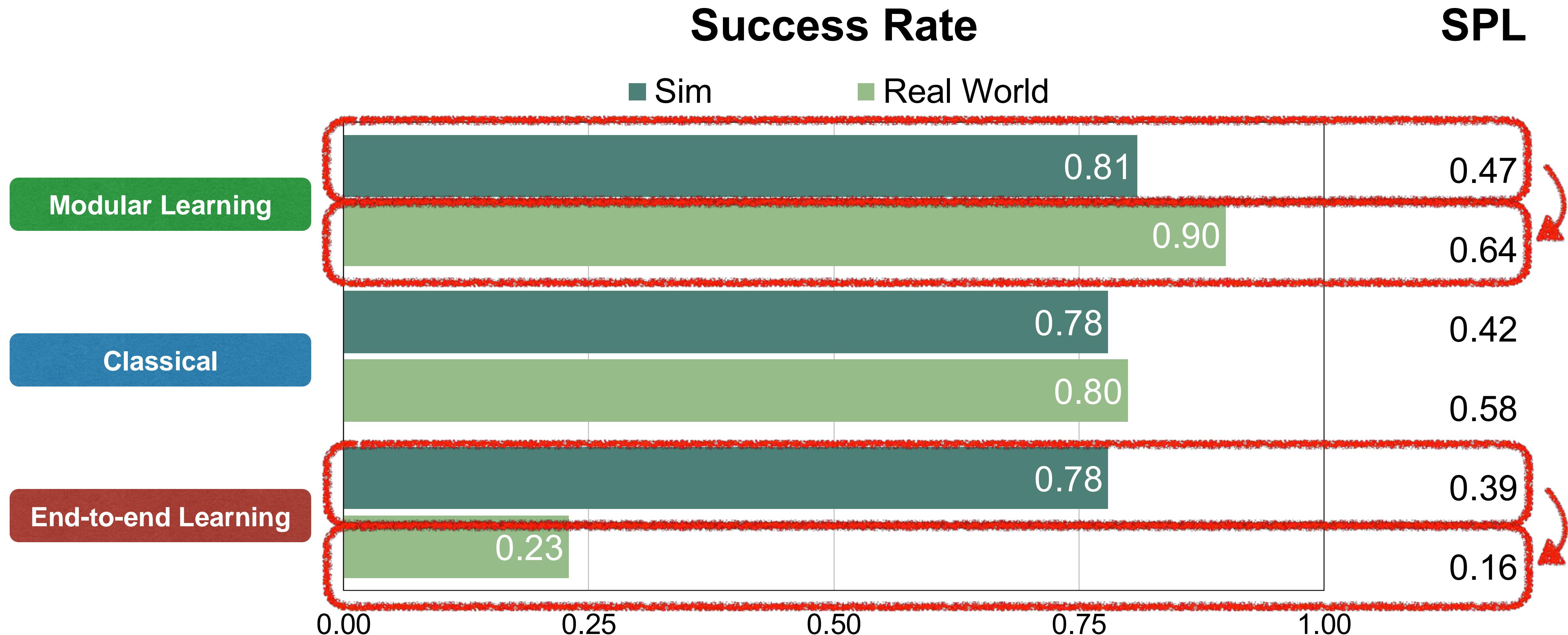
End-to-end fails to Transfer



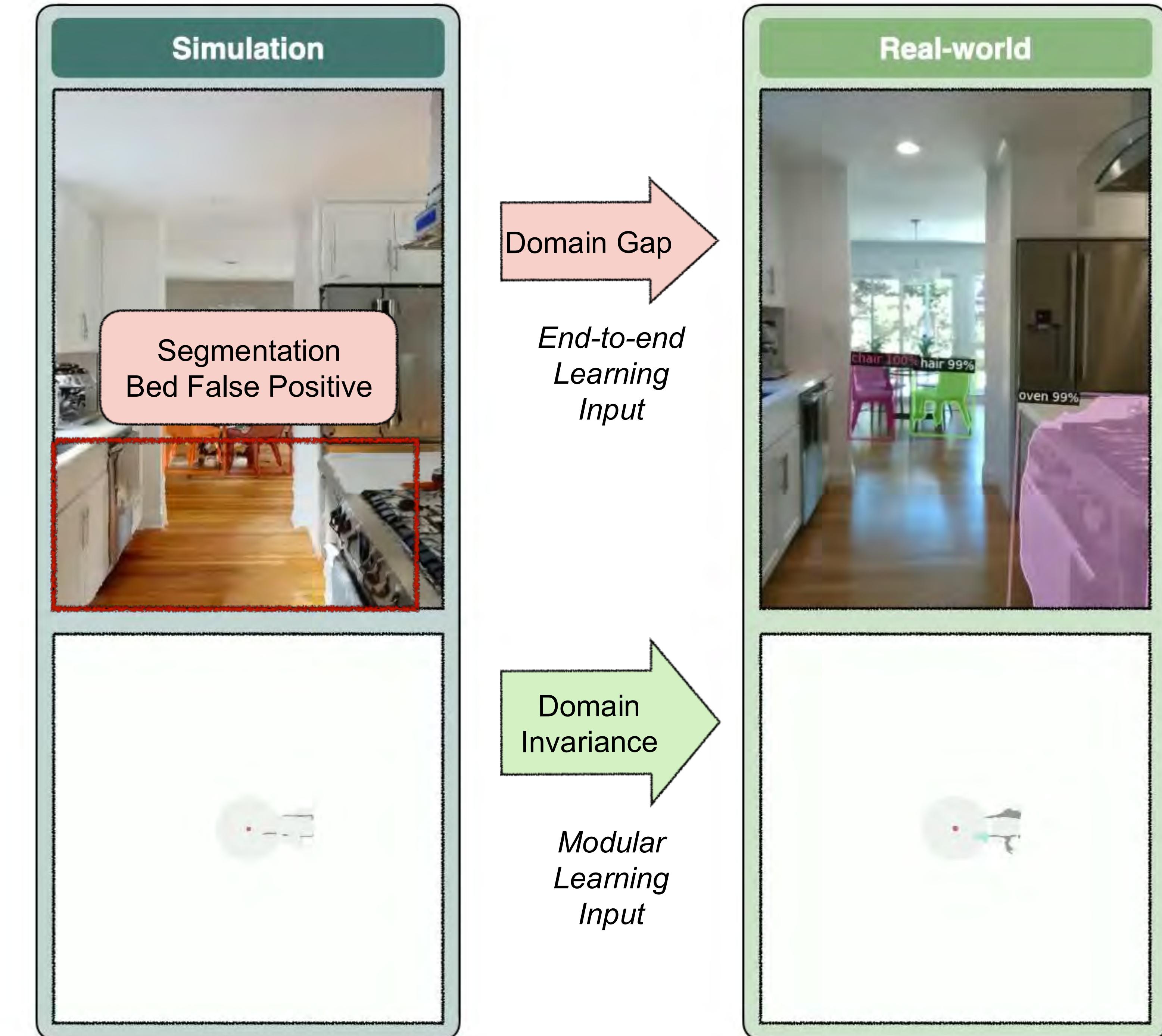
End-to-end Failures



Modular vs End-to-end Transfer

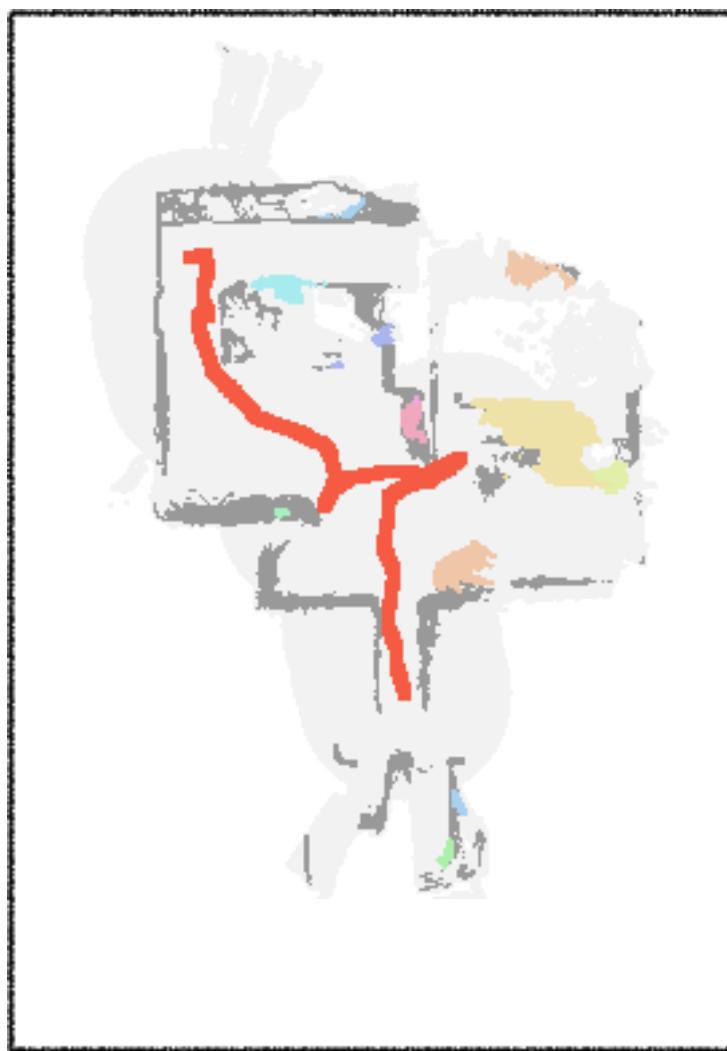


Simulation vs Reality



Real World

Predicted Semantic Map

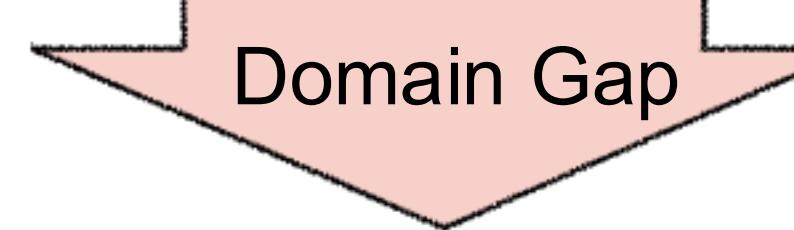


Domain Invariance

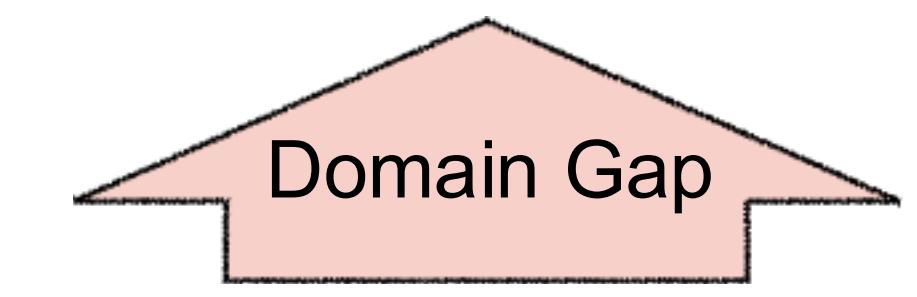
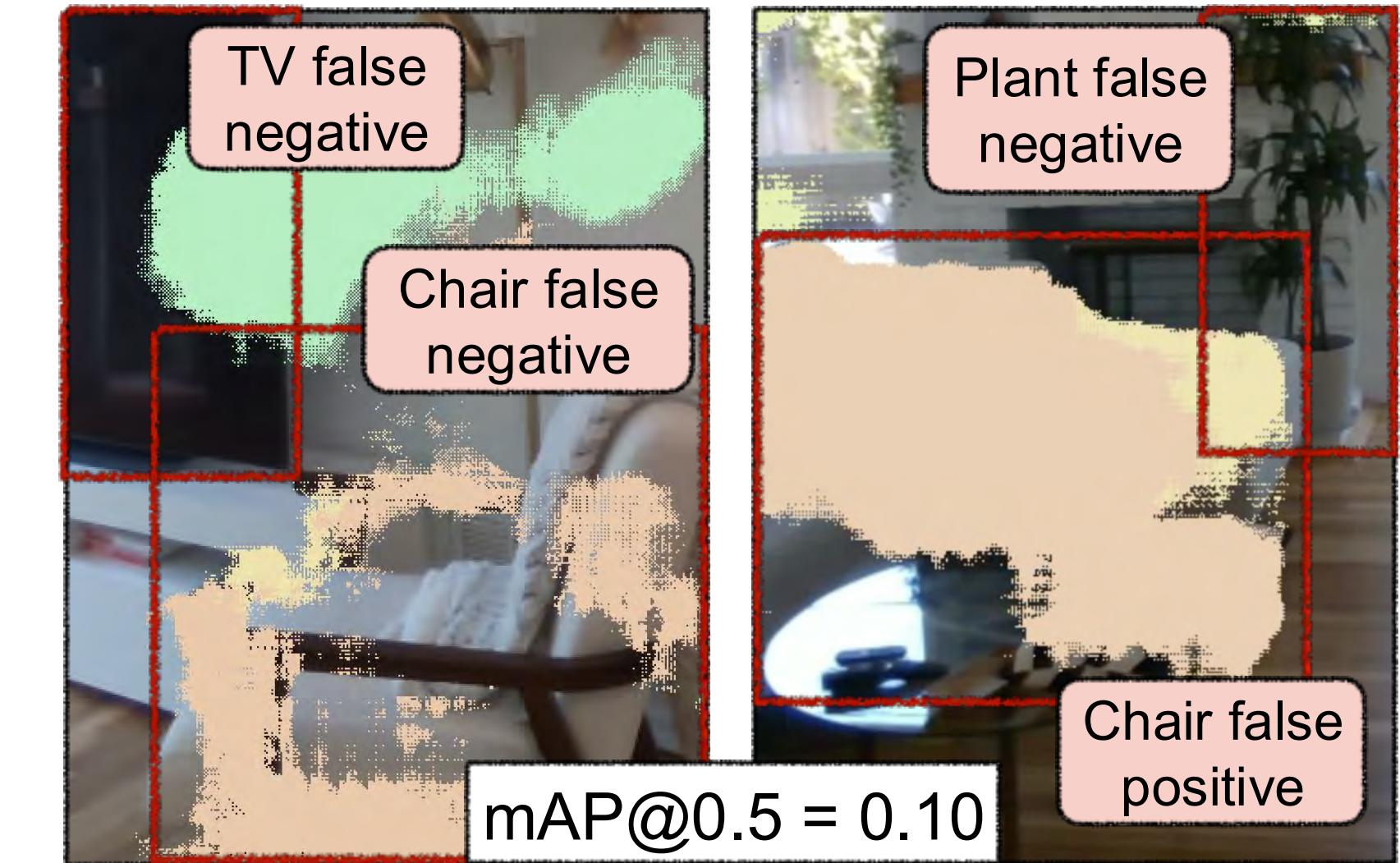
Simulation



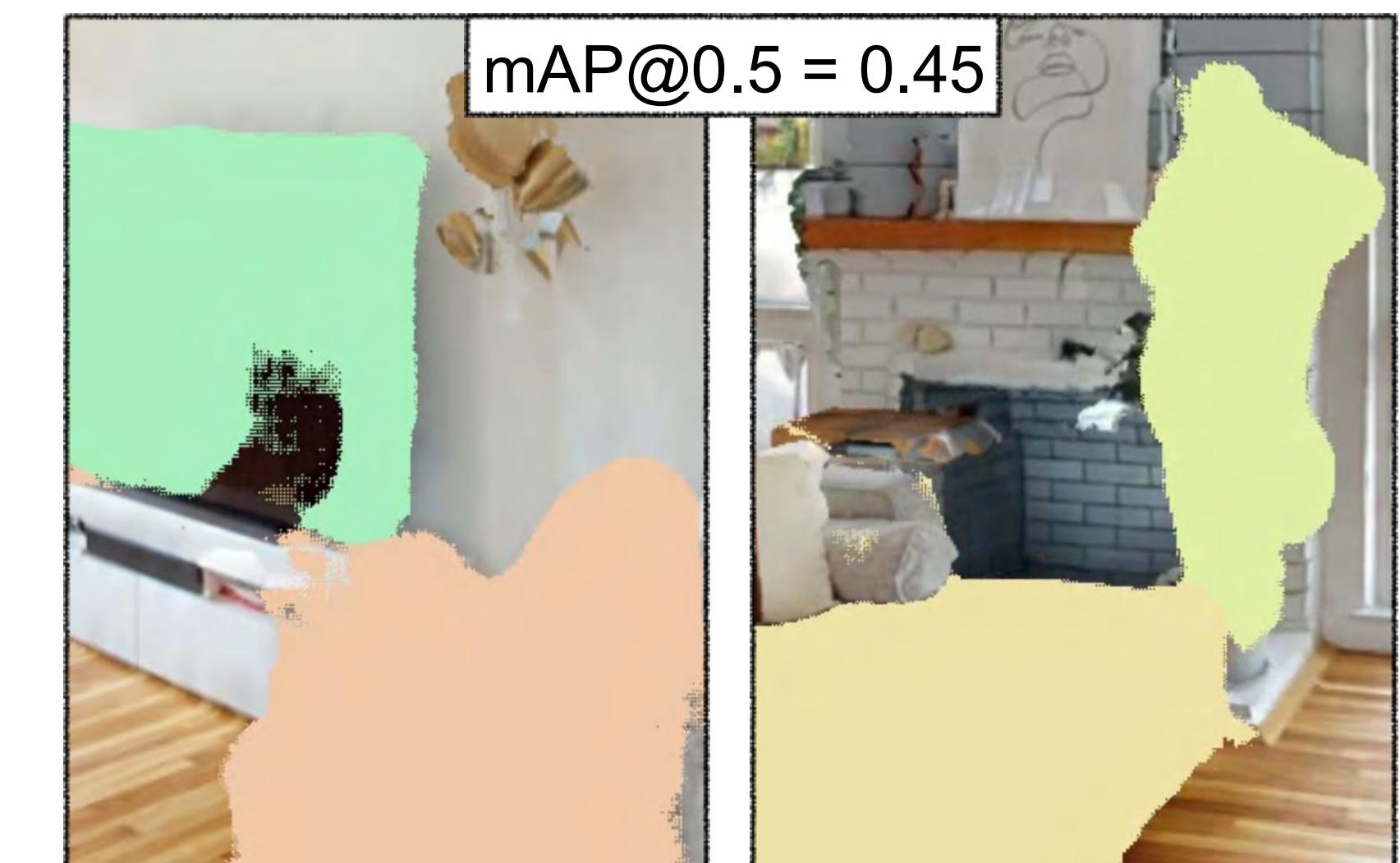
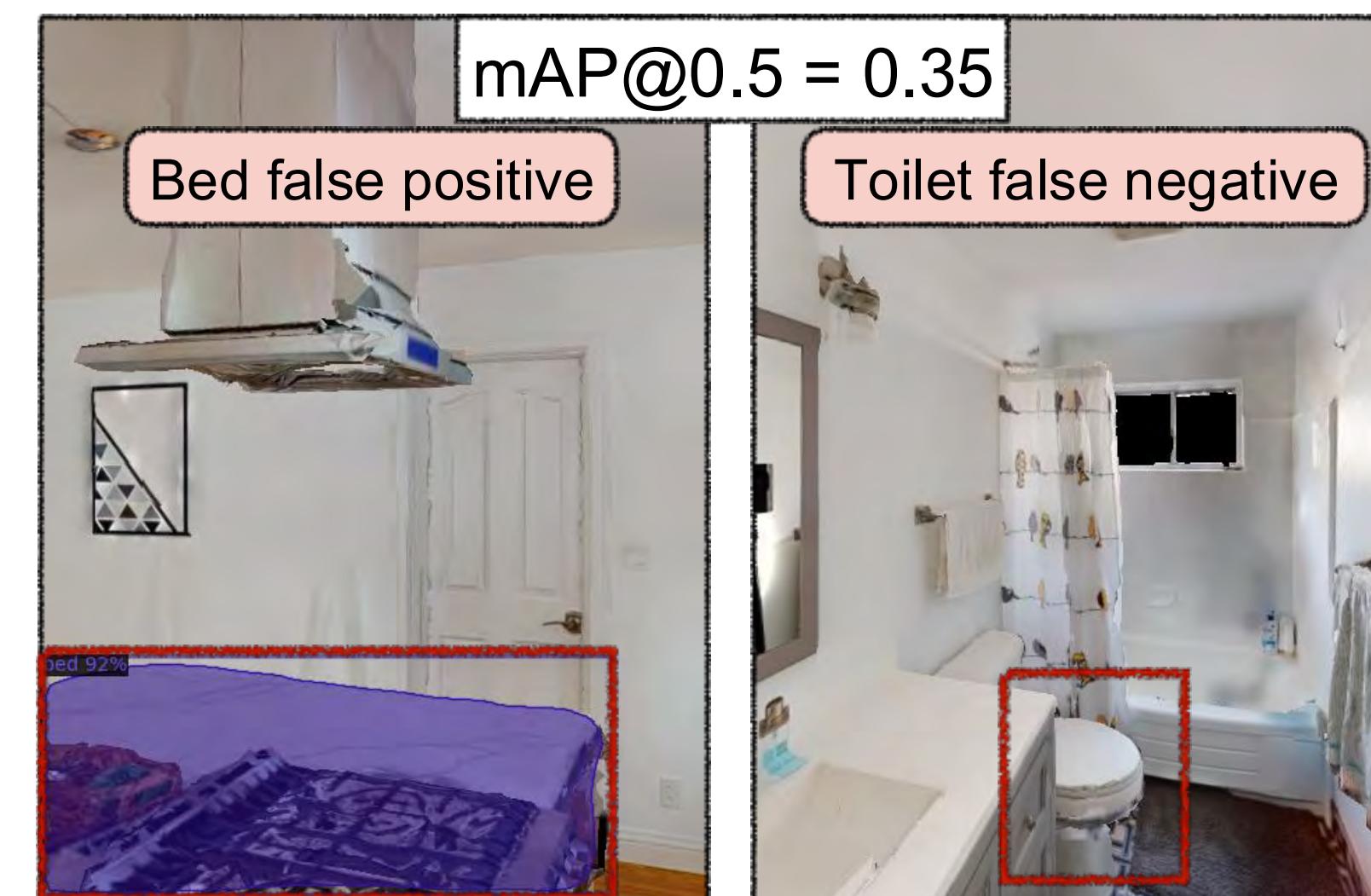
Segmentation Model Trained in Real World



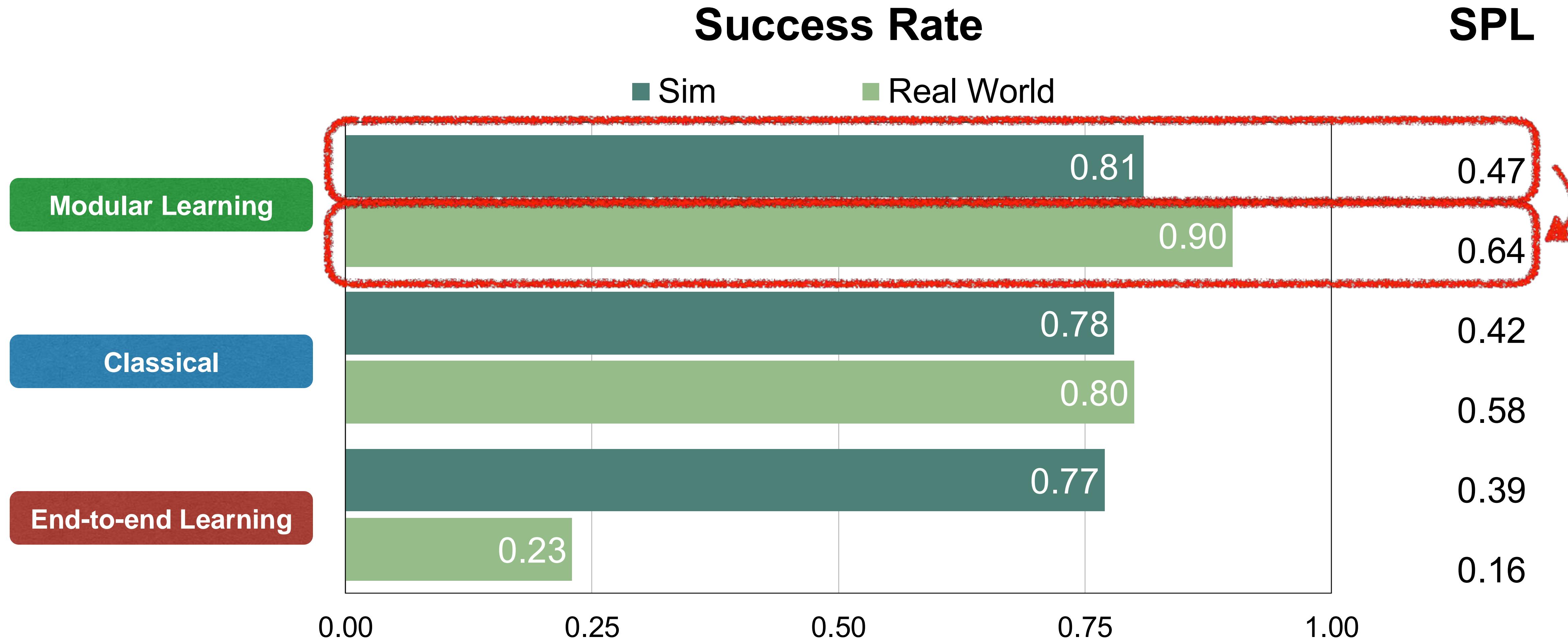
Segmentation Model Trained in Simulation



- 0: chair
- 1: couch
- 2: potted plant
- 3: bed
- 4: toilet
- 5: tv



Modular Learning Sim vs Real



Takeaways

For practitioners:

- Modular learning can reliably navigate to objects with 90% success

For researchers:

- Models relying on RGB images are hard to transfer from sim to real → *leverage modularity and abstraction in policies*
- Disconnect between sim and real error modes → *evaluate semantic navigation on real robots*



GOAT: GO to Any Thing

Matthew Chang*, Theophile Gervet*, Mukul Khanna*, Sriram Yenamandra*, Dhruv Shah, Tiffany Min, Kavit Shah, Chris Paxton, Saurabh Gupta, Dhruv Batra, Rozbeh Mottaghi, Jitendra Malik*, Devendra Singh Chaplot*

Third-person view



1. GOAT Problem

GOAT System Architecture

Results

Applications

Pick & Place

Social Navigation

Platform Agnostic

Unknown Environment
Explore

Perception
*Detect and Localize
Objects*

Lifelong Memory
*Remember Object
Locations*

Control
*Navigate to / Pick & Place
Objects*



Multimodal:

Reach Any Object Specified in Any Way

Image



Language

*Find **the fruit basket** on the kitchen counter*

Category

Bring me a CUP

Lifelong:

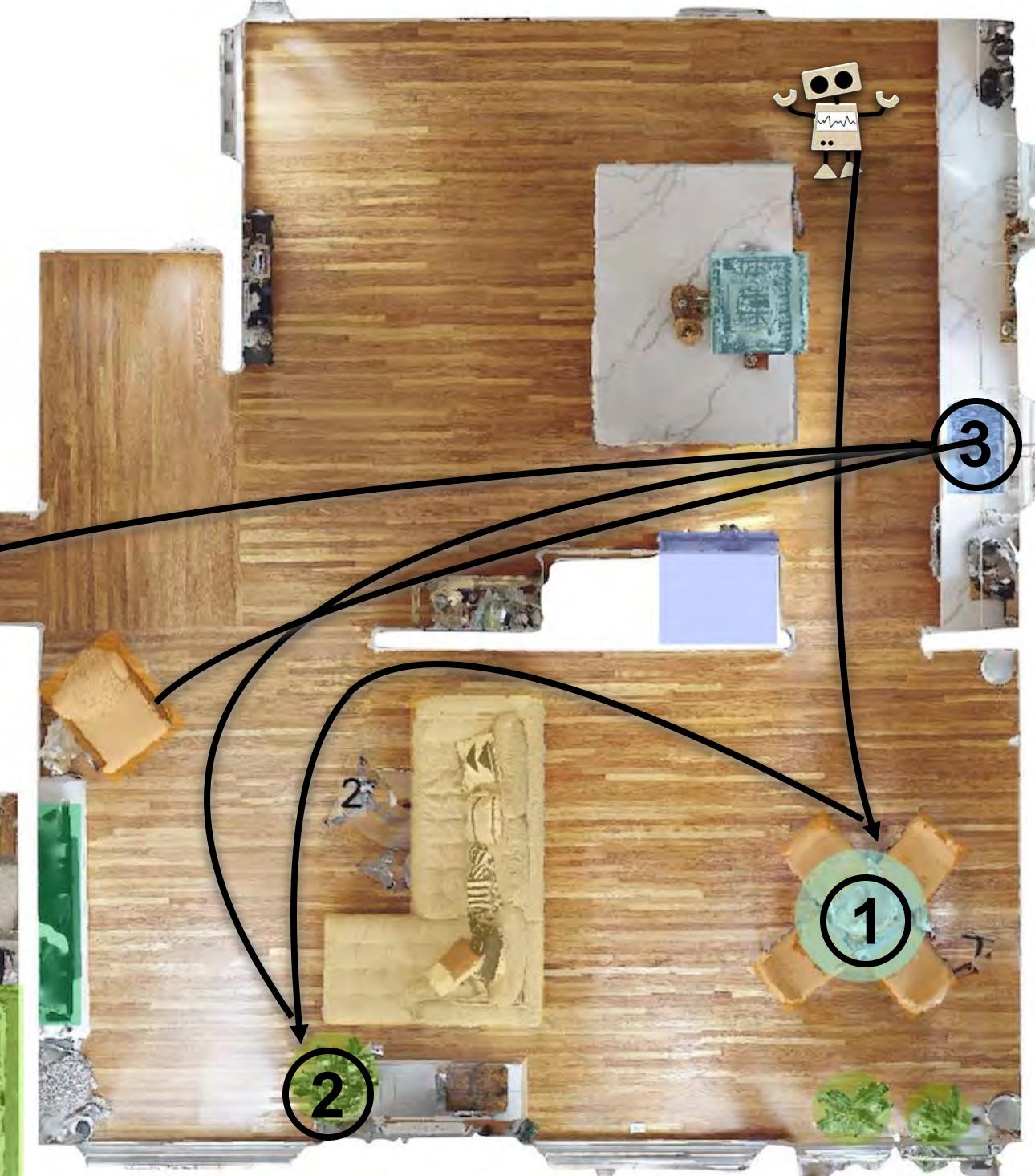
Remember Object Locations



① Go to the potted plant next to the couch

② Go to a SINK

③ Go to the black and white striped bed



0 1 2 3m
10ft

1. GOAT Problem

GOAT System Architecture

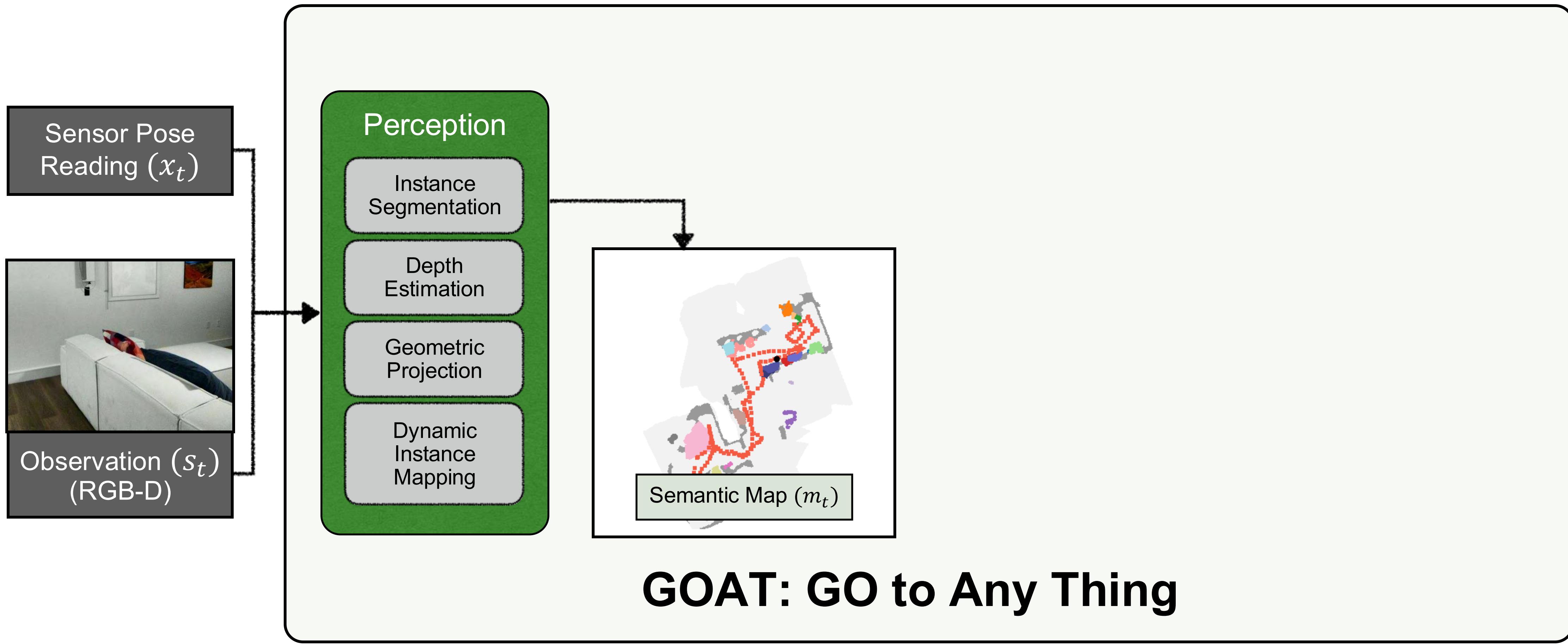
Results

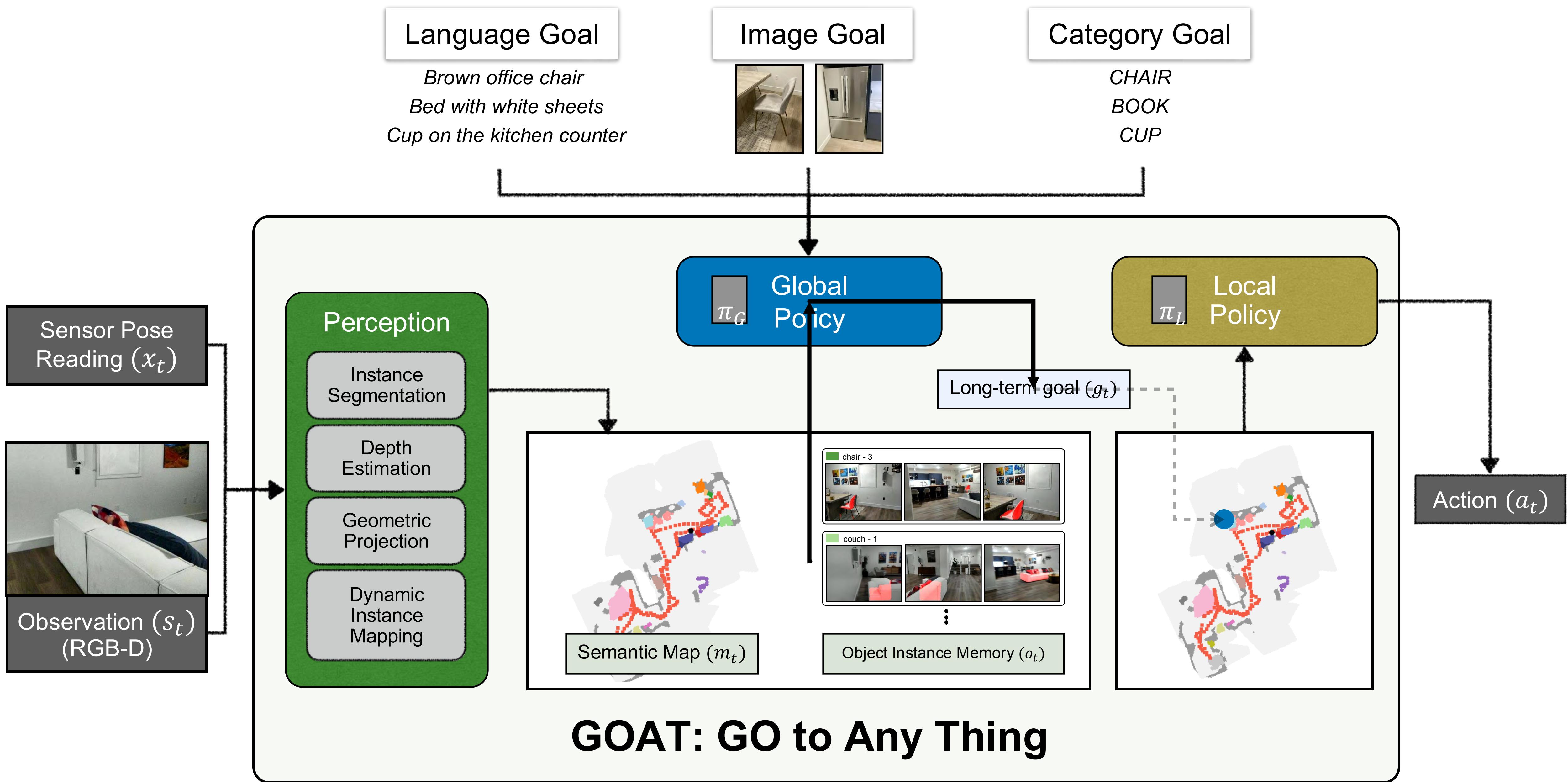
Applications

Pick & Place

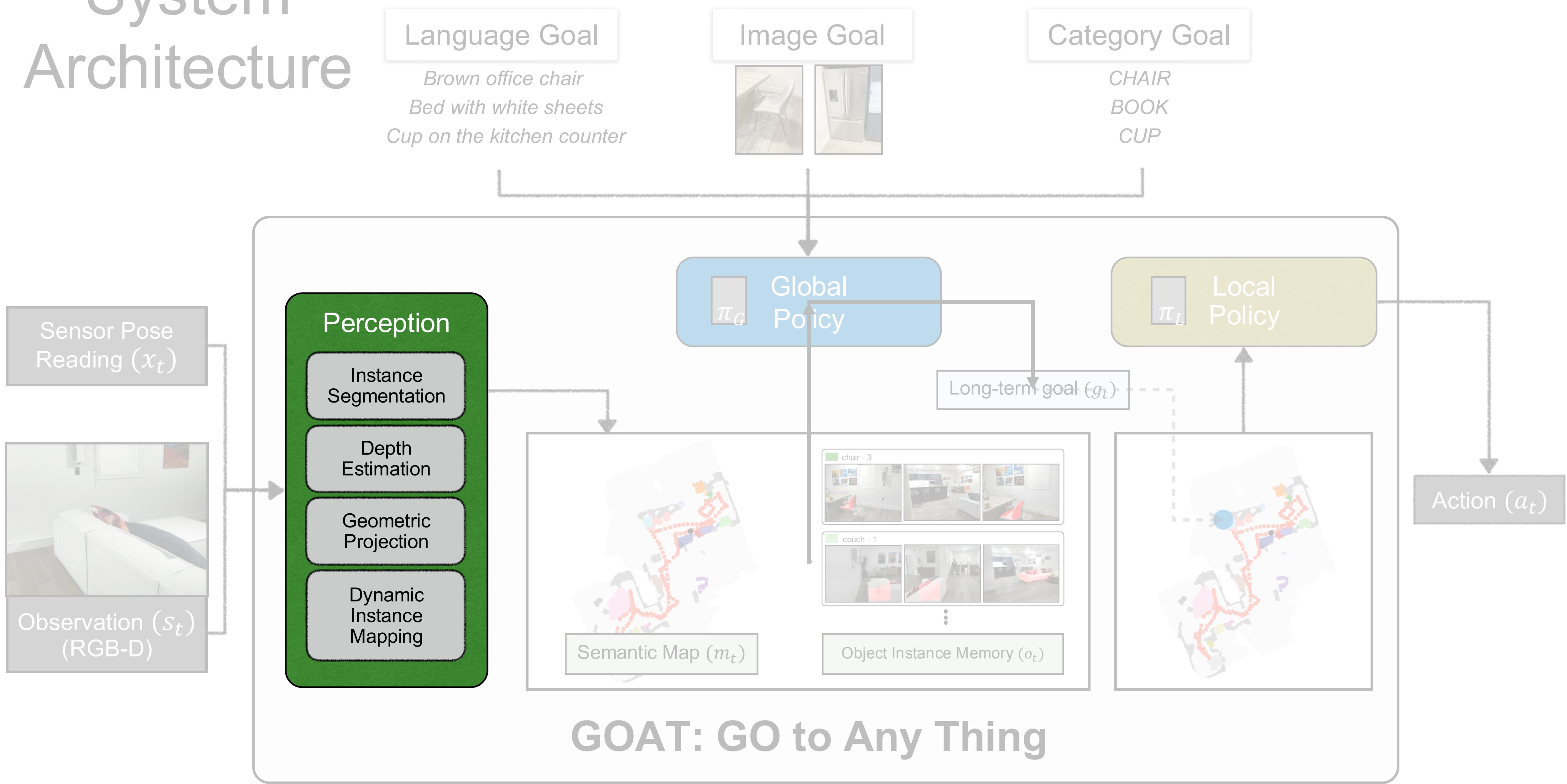
Social Navigation

Platform Agnostic

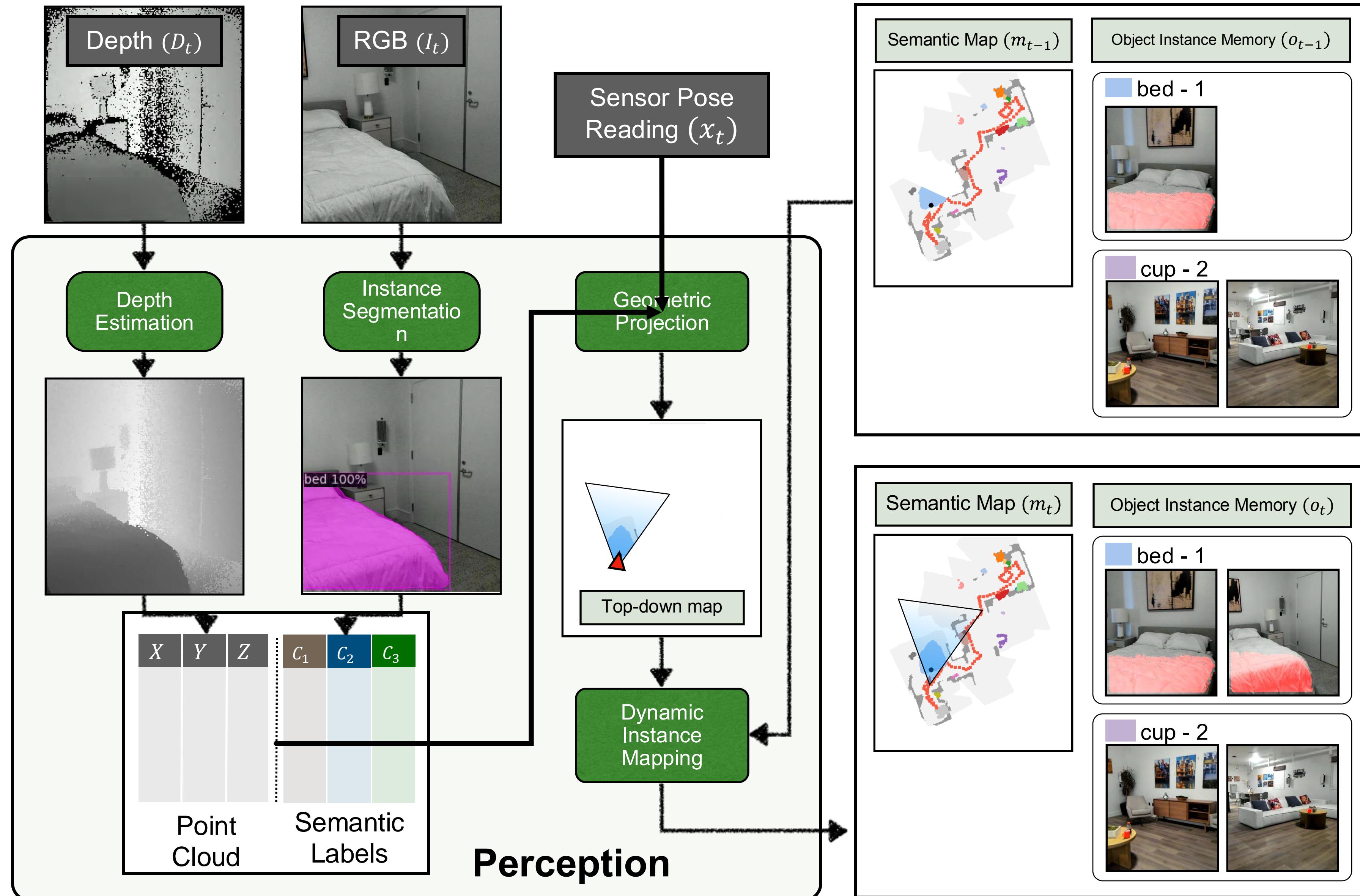




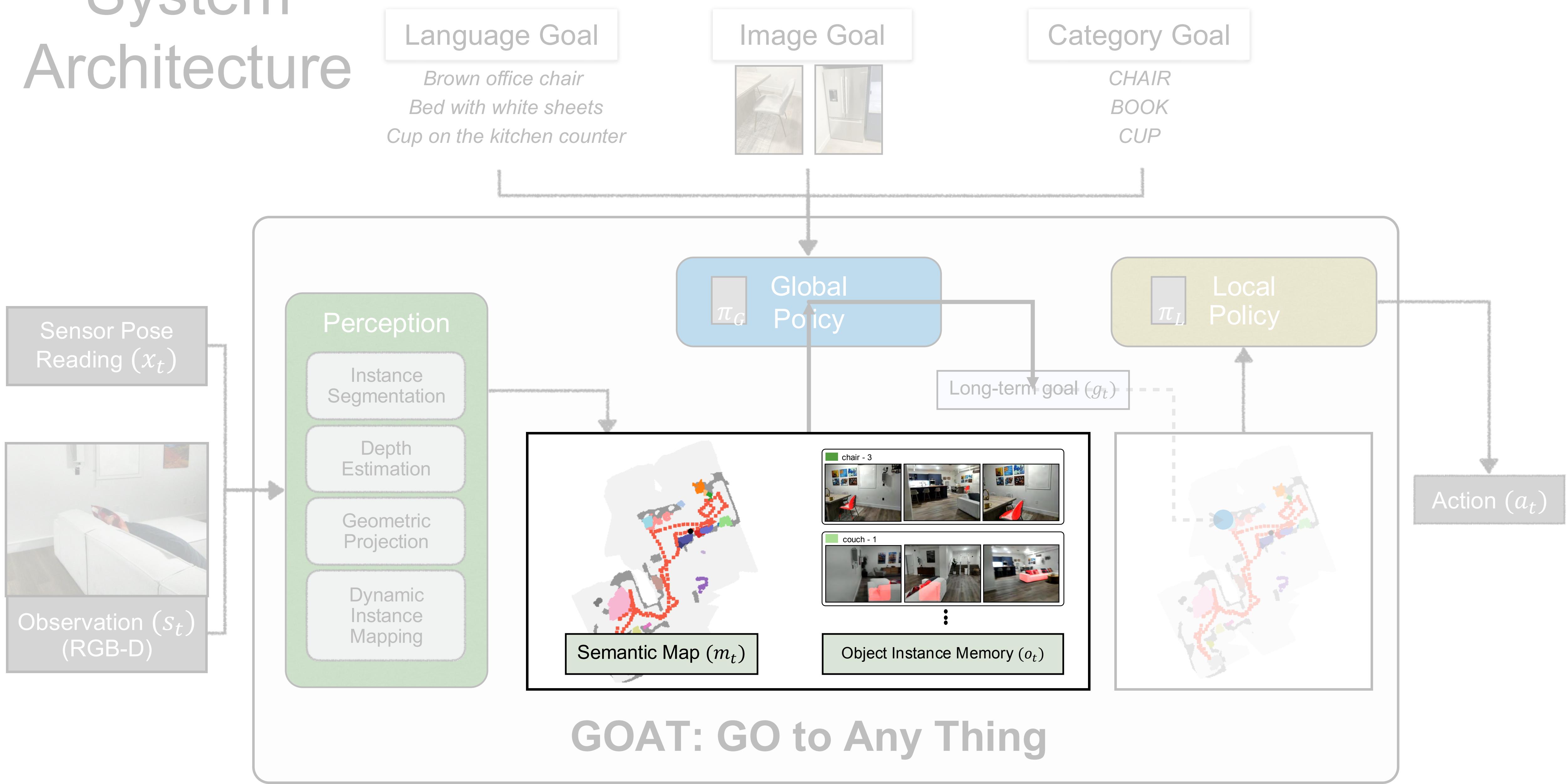
System Architecture



Perception System

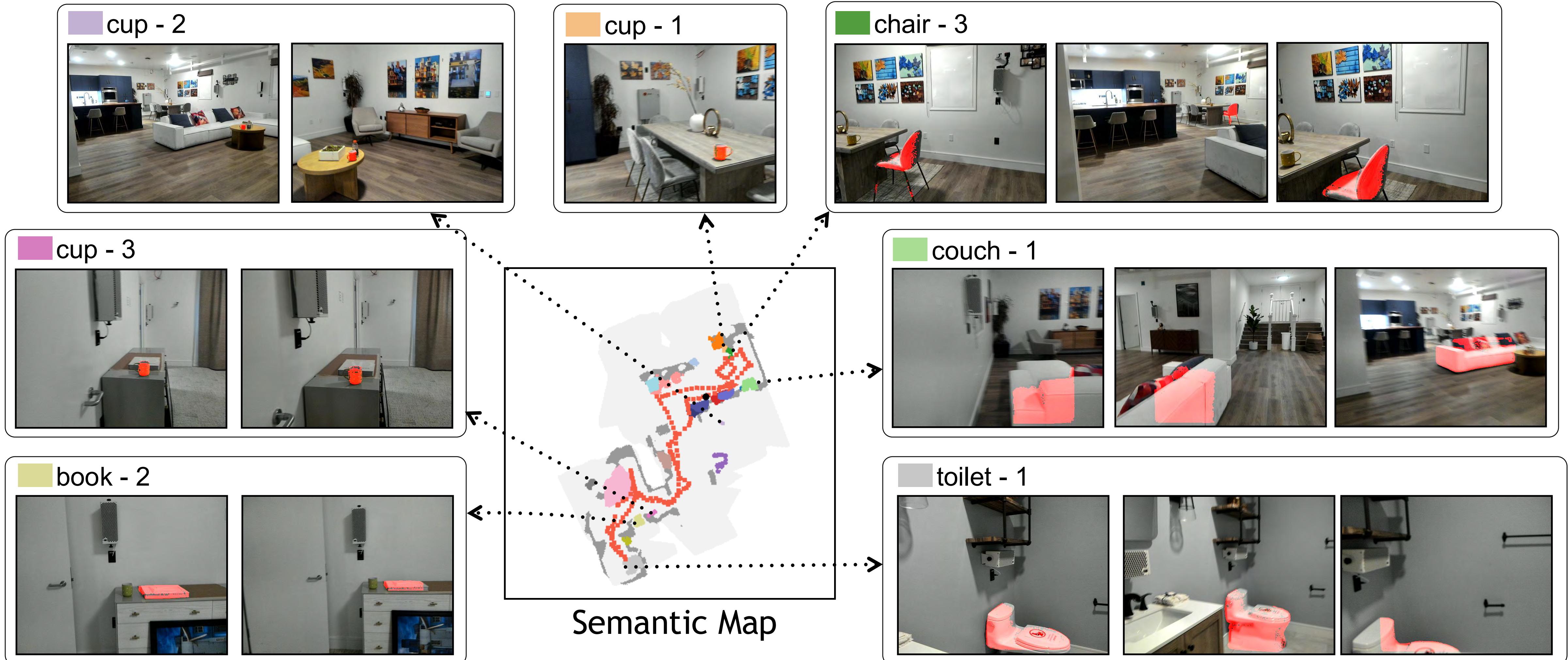


System Architecture

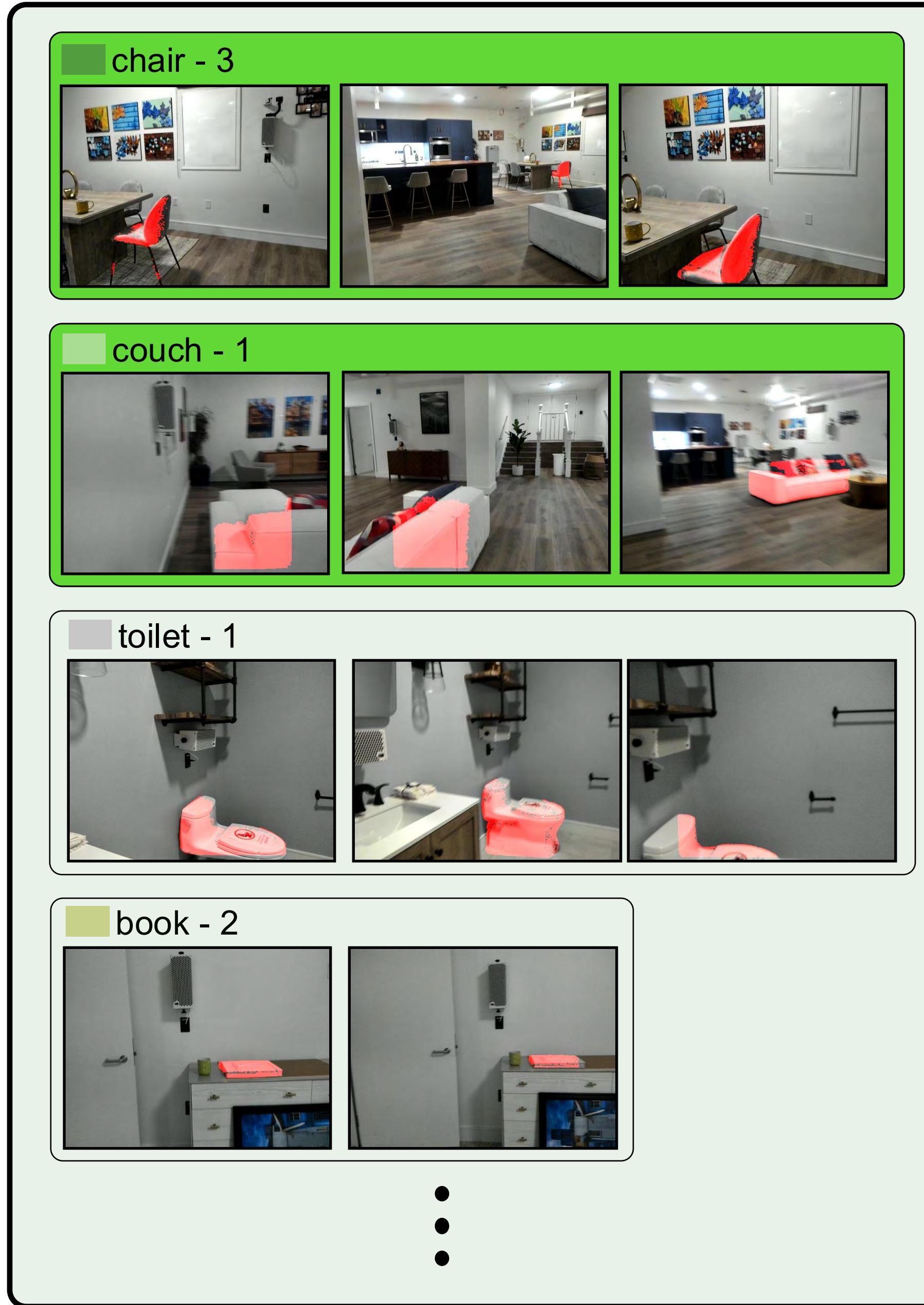


GOAT Memory Representation

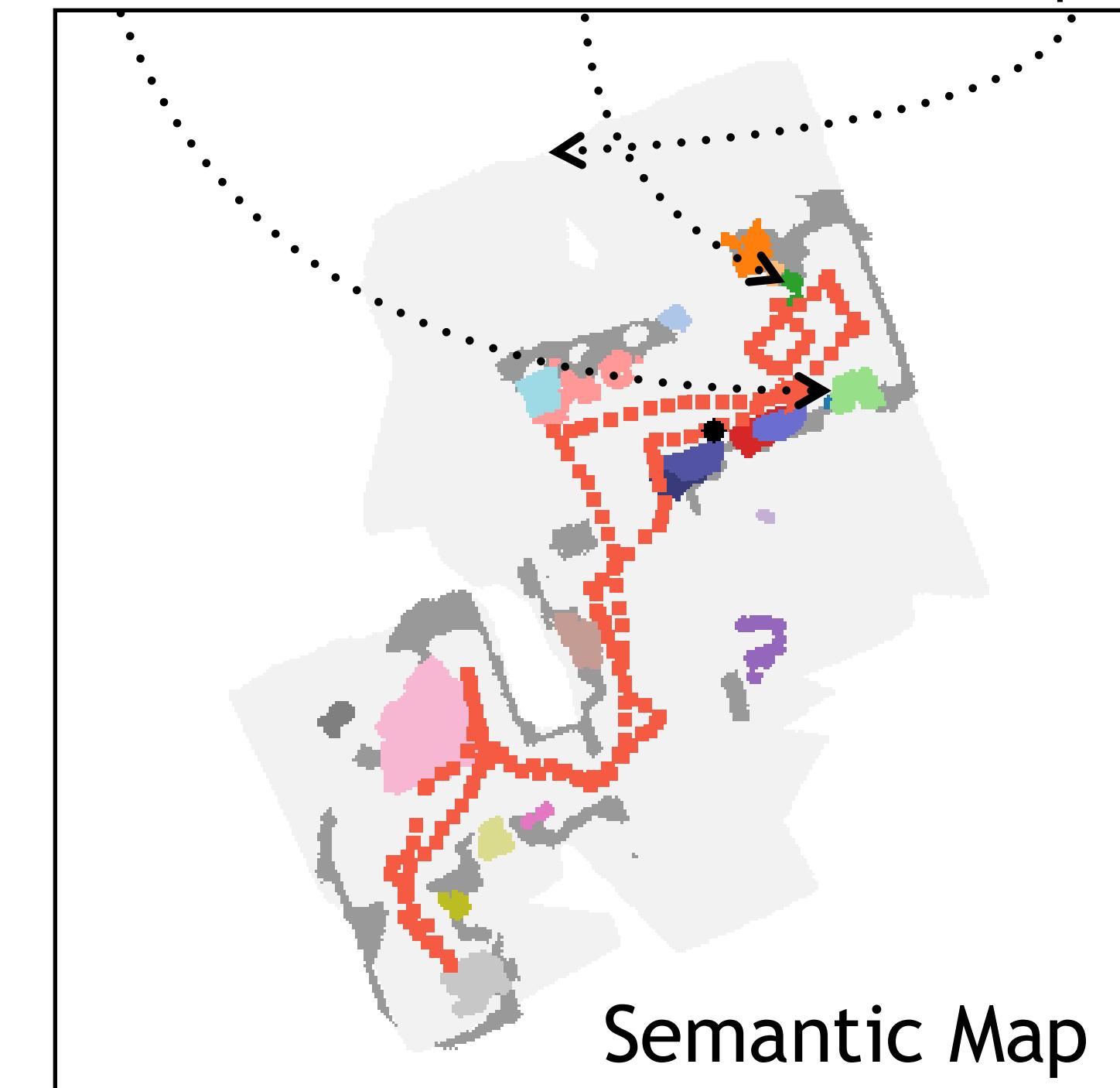
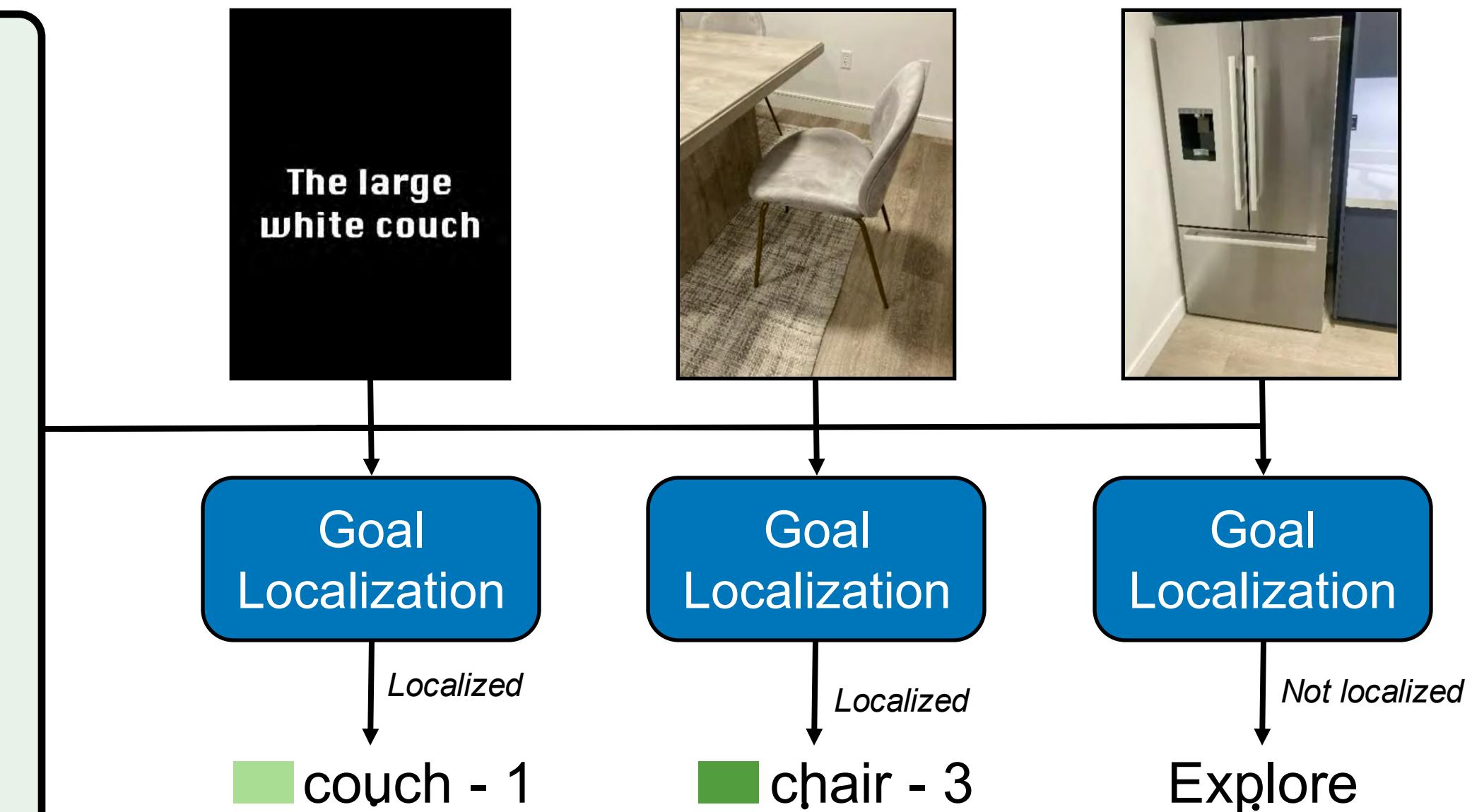
Semantic map with associated Object Instance memory



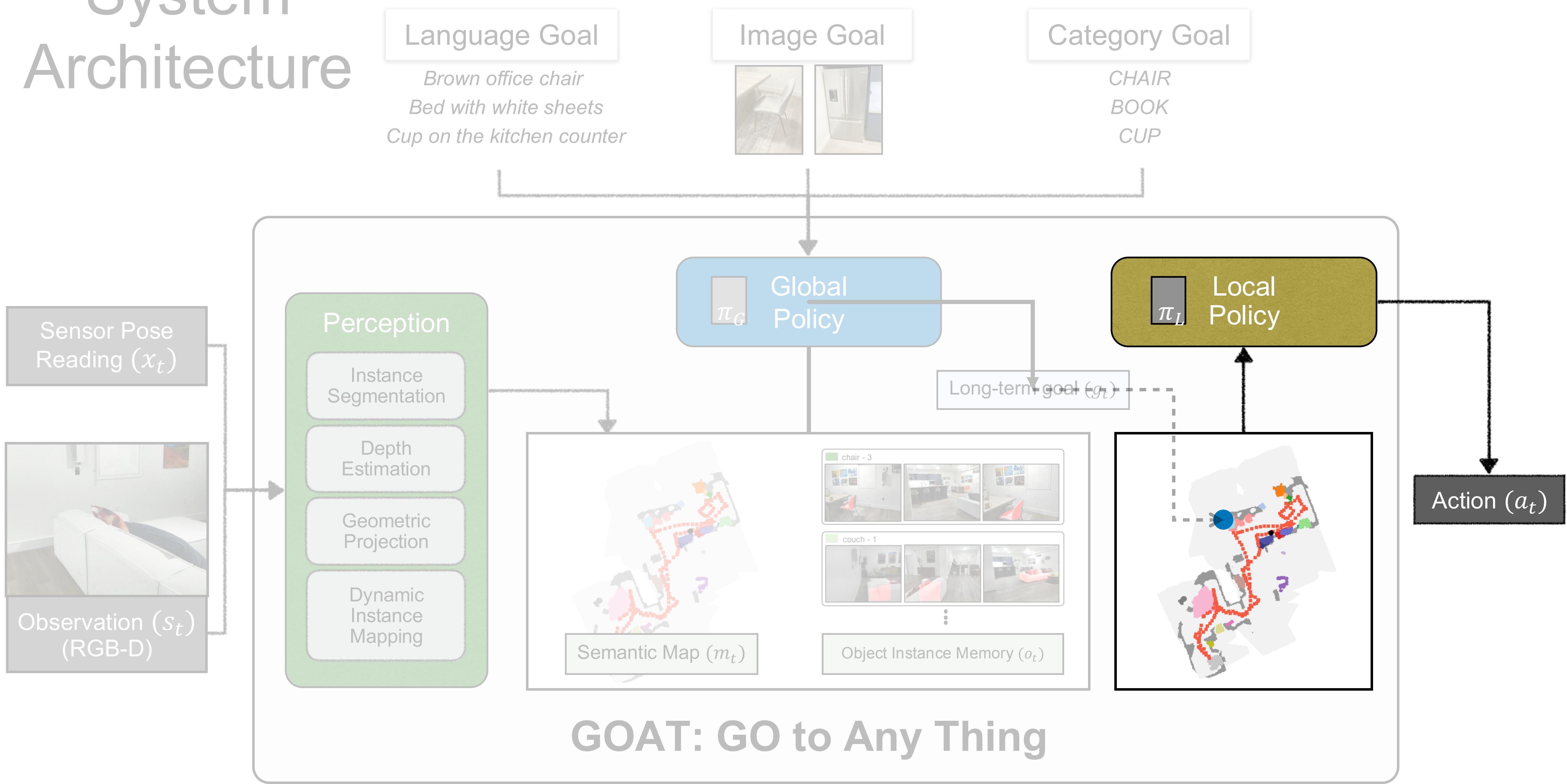
Object Instance Memory



Goals



System Architecture



1. GOAT Problem

GOAT System Architecture

Results

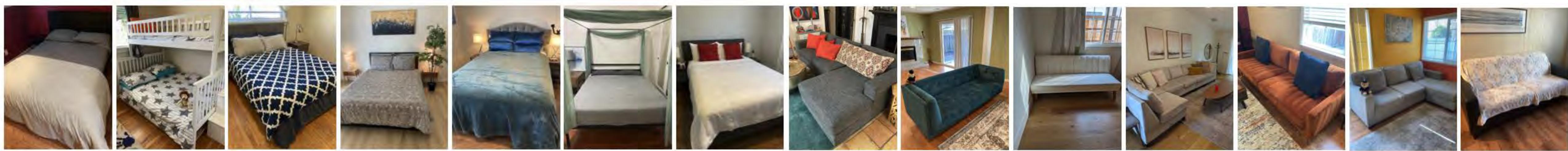
Applications

Pick & Place

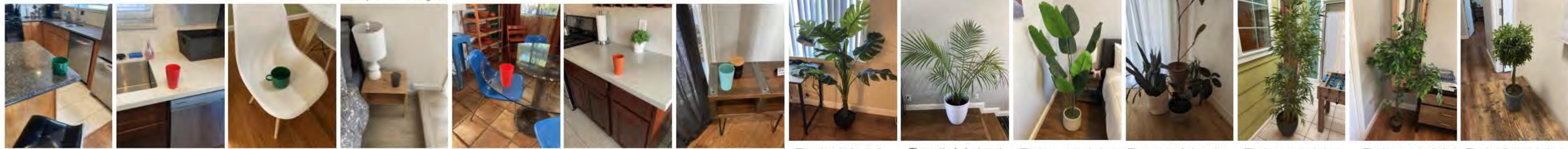
Social Navigation

Platform Agnostic

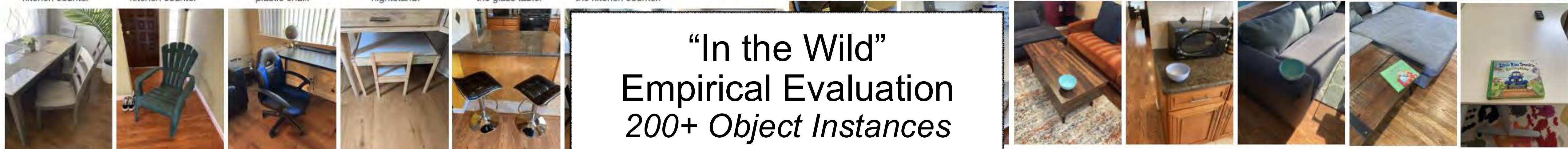
"In the Wild"
Empirical Evaluation
9 Unseen Homes
4 Methods
10 Trajectories per Home
5-10 Goals per Trajectory
~90h of Experiments



The bed with the white blanket pulled back halfway and grey sheets.
The bunk bed with stars on the blanket.
The bed with blue and white sheets.
The with the blue and yellow painting above it and plant on the right.
The bed with the blue blanket and blue pillows.
The bed with grey sheets and light green trim.
The bed with the white blanket and red pillows
The large grey living room couch with many pillows.
The green couch.
The white rectangular couch with no pillows.
The large grey living room couch with many pillows.
The light brown couch with blue pillows.
The large grey couch in front of the yellow wall.
The couch covered in a white blanket.



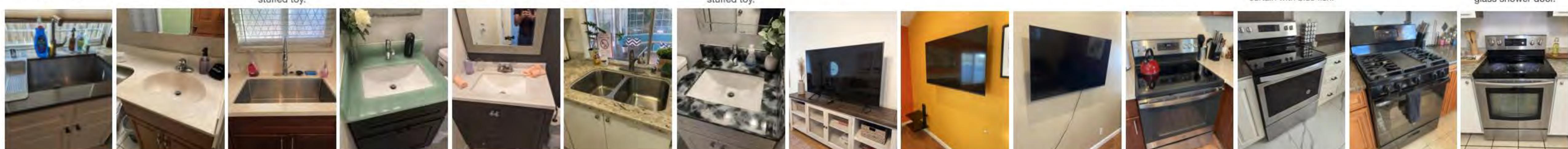
The green cup on the kitchen counter
The red cup on the kitchen counter
The green mug on the plastic chair.
The grey cup on the nightstand.
The red cup on top of the glass table.
The red cup on the kitchen counter.
The light blue cup.
The plant in front of the window.
The potted plant next to the stairs.
The large potted plant in the bedroom.
The group of plants in front of the curtain.
The large potted plant next to the foosball table.
The large potted plant in front of the mirror.
The small potted plant on the hallway table



The grey dining table chair.
The green lawn chair.
The black office chair.
The small wooden chair at the desk.
The black leather chairs in the kitchen.
The bar-height chair at the kitchen island.
a wooden seat.
wood desk.
the kitchen counter
The light blue bowl on the living room coffee table.
The lavender bowl on the kitchen counter.
The light blue bowl on the couch.
The green cover book on the coffee table.
The book on the desk. It has a car on the cover.



The beige teddy bear.
The stuffed lion toy.
The green dinosaur stuffed toy.
The stuffed lion.
The beige teddy bear.
The stuffed lion toy.
The green dinosaur stuffed toy.
The refrigerator.
The refrigerator.
The refrigerator.
The refrigerator.
The toilet next to the shower curtain with blue fish.
The toilet.
The toilet next to the glass shower door.



The kitchen sink.
The bathroom sink with marble top.
The kitchen sink.
The bathroom sink with green counter.
The bathroom sink.
The kitchen sink.
The bathroom sink with black and white counter.
The television.
The television mounted on a yellow wall.
The television mounted on a white wall.
The oven.
The oven.
The oven.
The oven.

"In the Wild"

Empirical Evaluation

200+ Object Instances

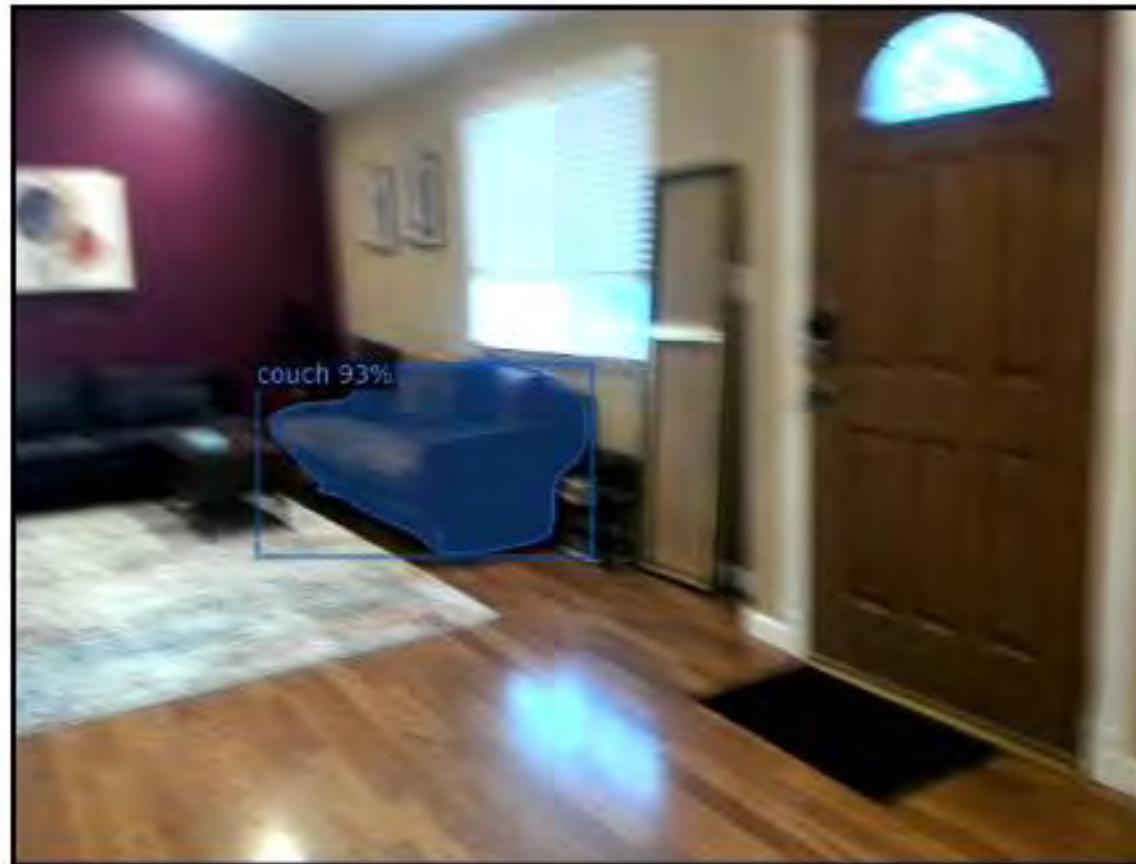
Observation



Third-person view



Observation



Third-person view



Semantic Map



couch - 1	chair - 1	chair - 2	couch - 2	chair - 3	chair - 4	chair - 5
refrigerator - 1	chair - 6	dining table - 1	potted plant - 1	potted plant - 2	sink - 1	oven - 1
bottle - 1	bottle - 2	chair - 7	tv - 1	chair - 8	teddy bear - 1	chair - 9
cup - 1						

Third-person view

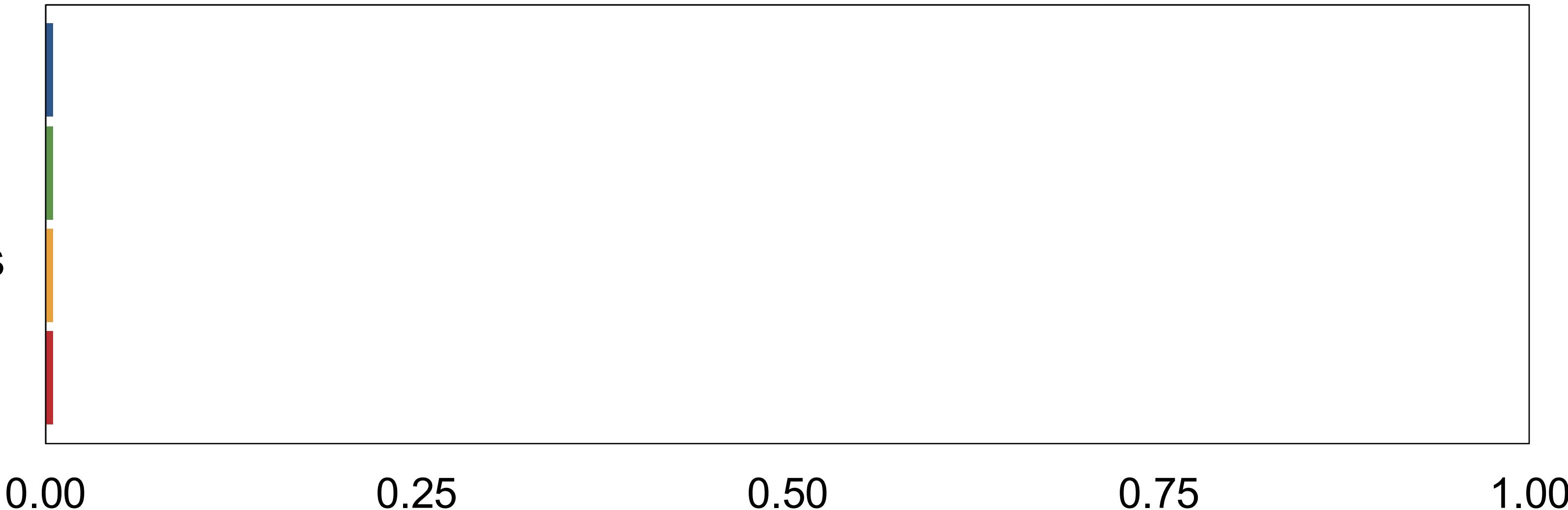


Results

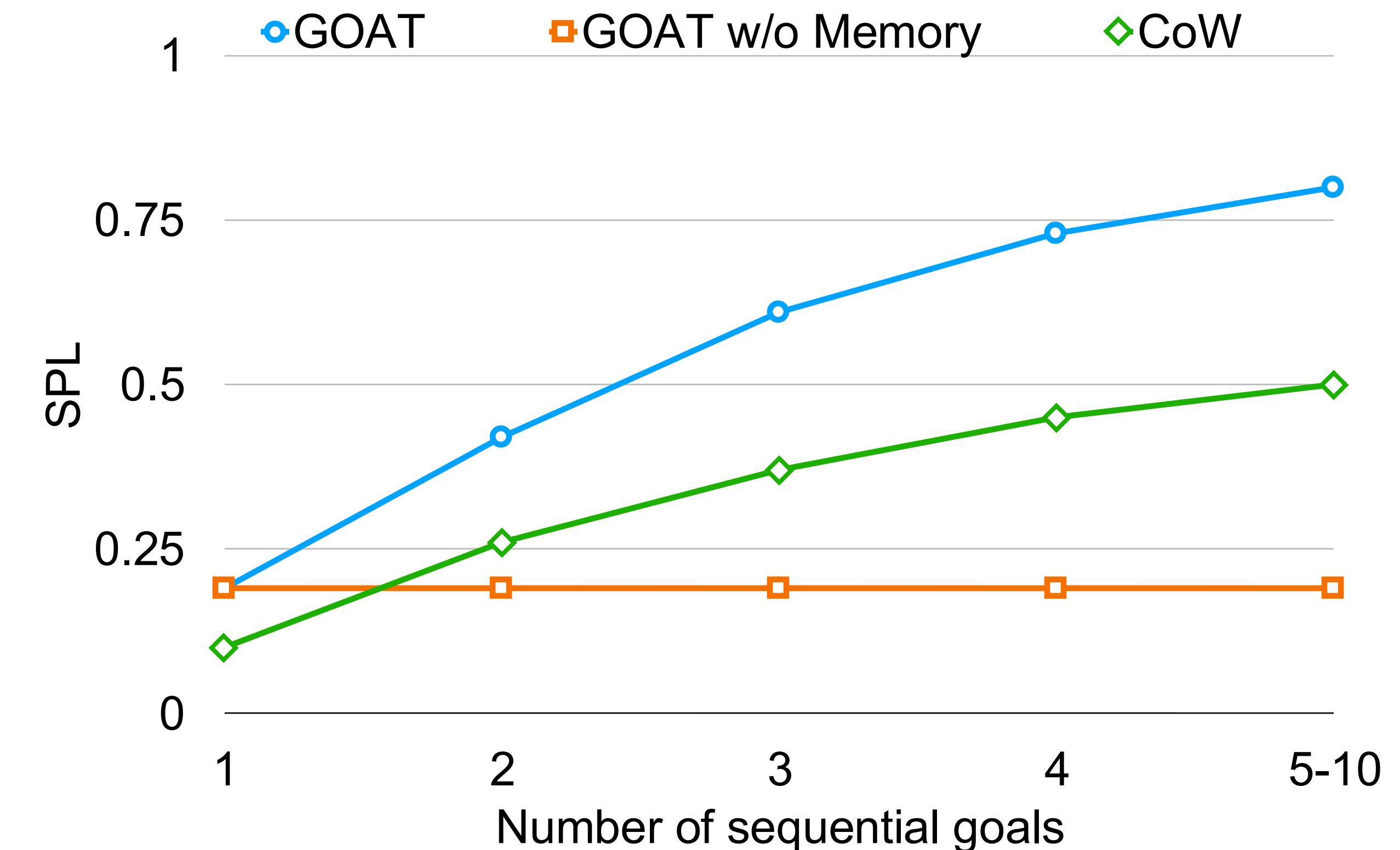
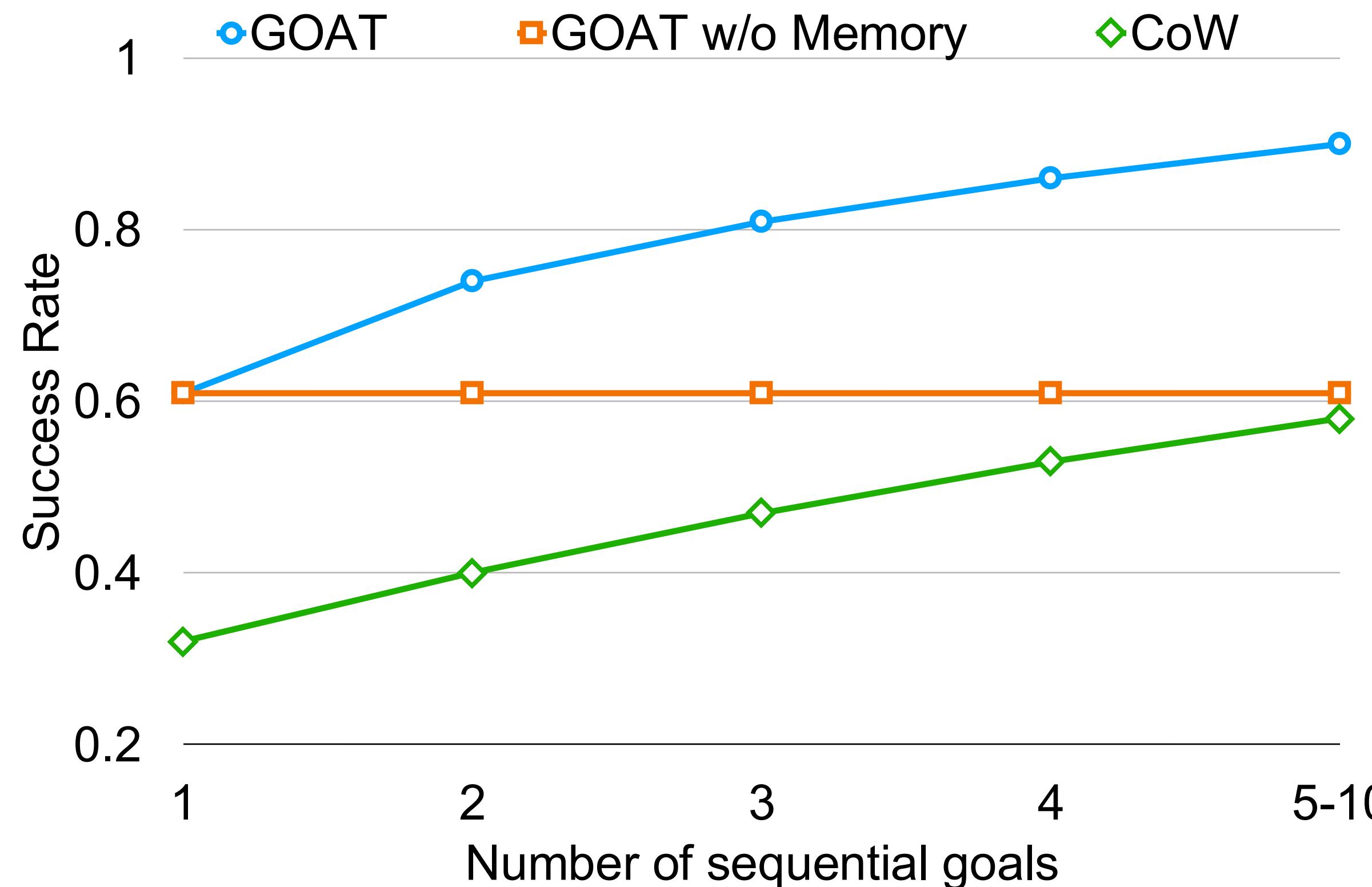
Success Rate

SPL

- GOAT
- GOAT w/o Memory
- GOAT w/o Instances
- CLIP on Wheels



Performance Across Episode



Goal:



Baselines

Ours

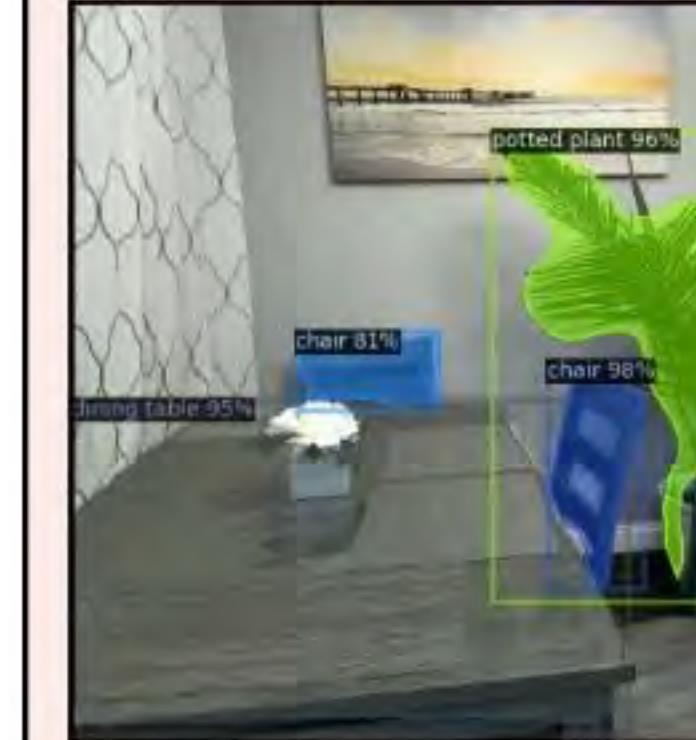
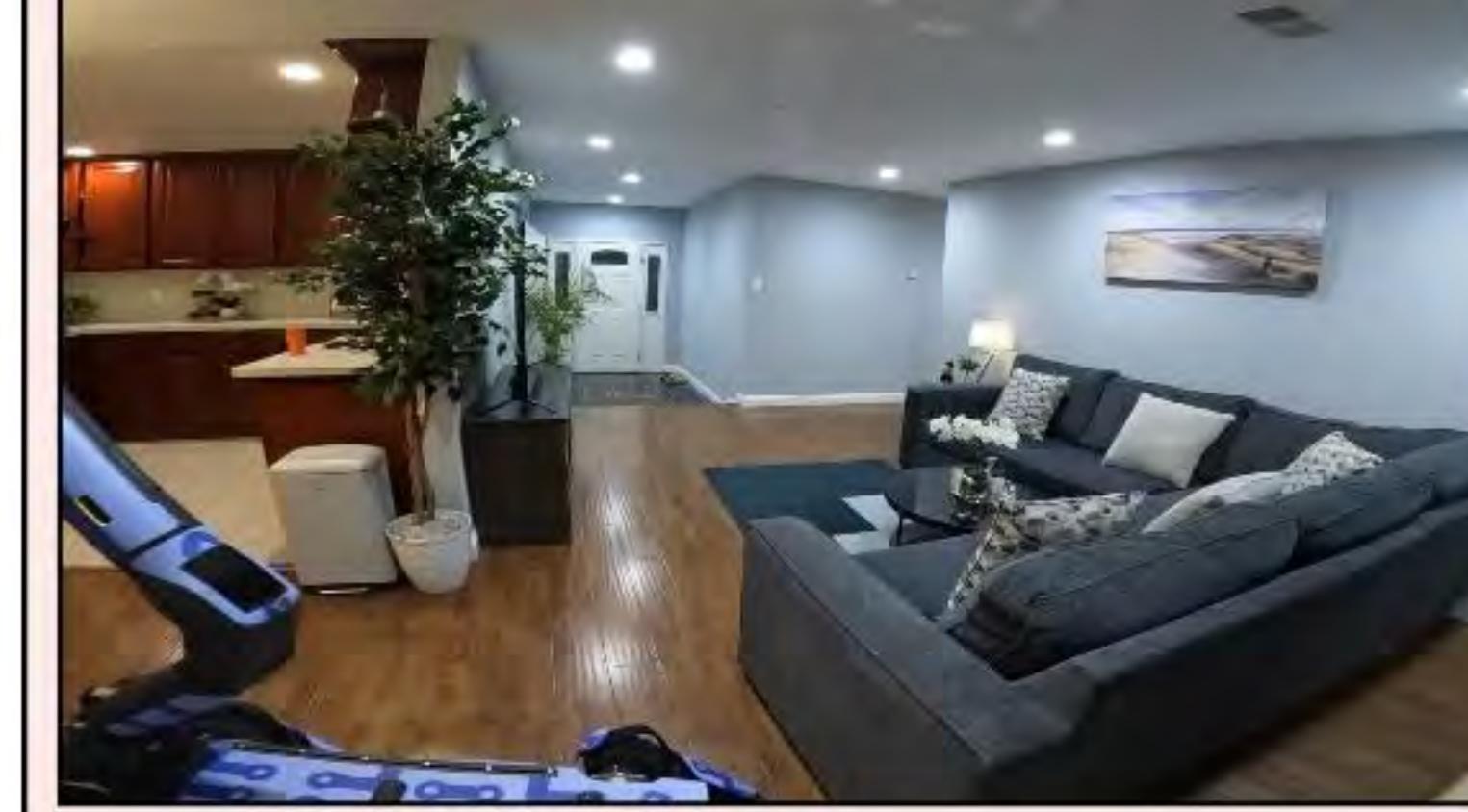


Observation



Instance Map

No Memory

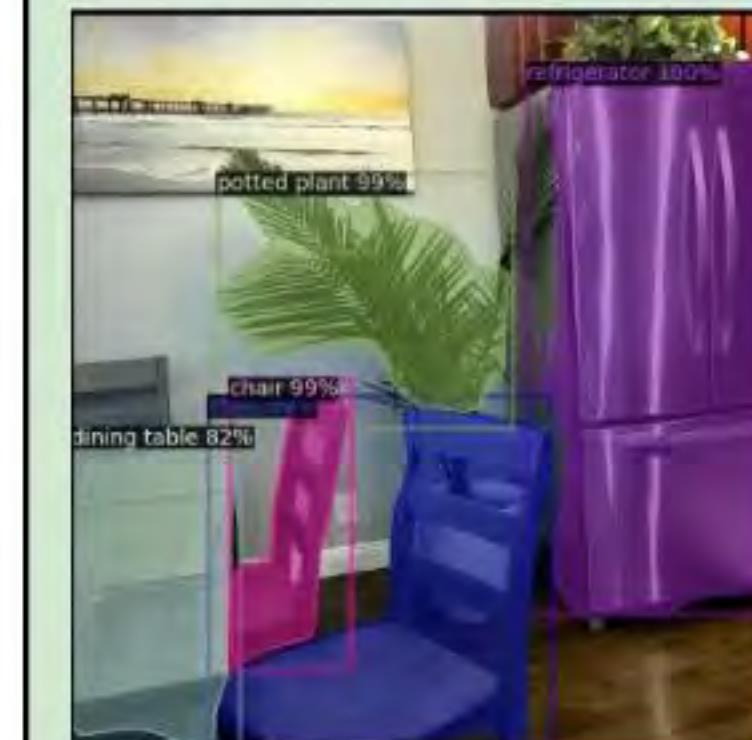


Observation



Instance Map

CoW



Observation



Instance Map

Success: 6/6 SPL: 0.78

Success: 4/6 SPL: 0.40

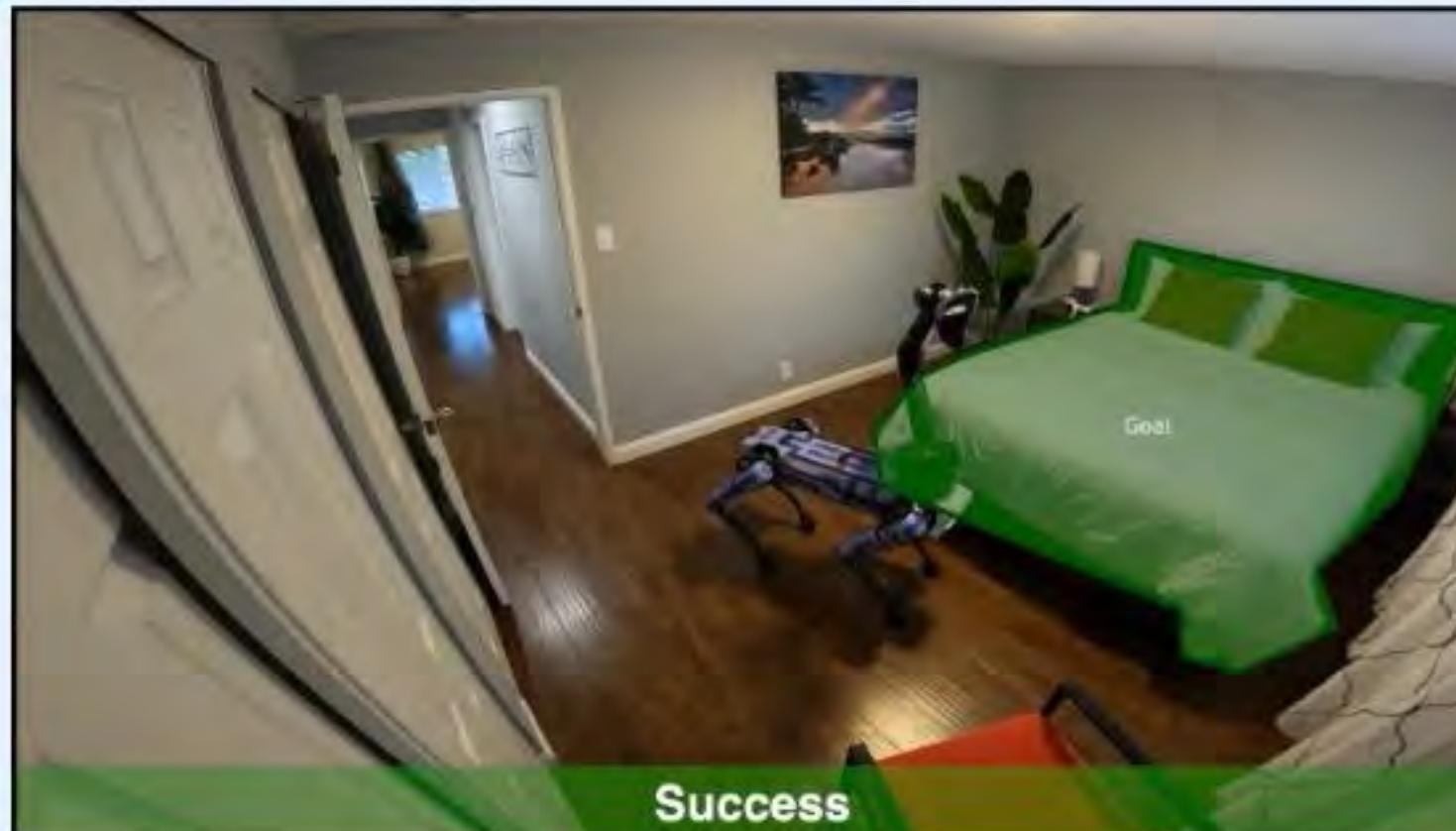
Success: 1/6 SPL: 0.16

Goal:

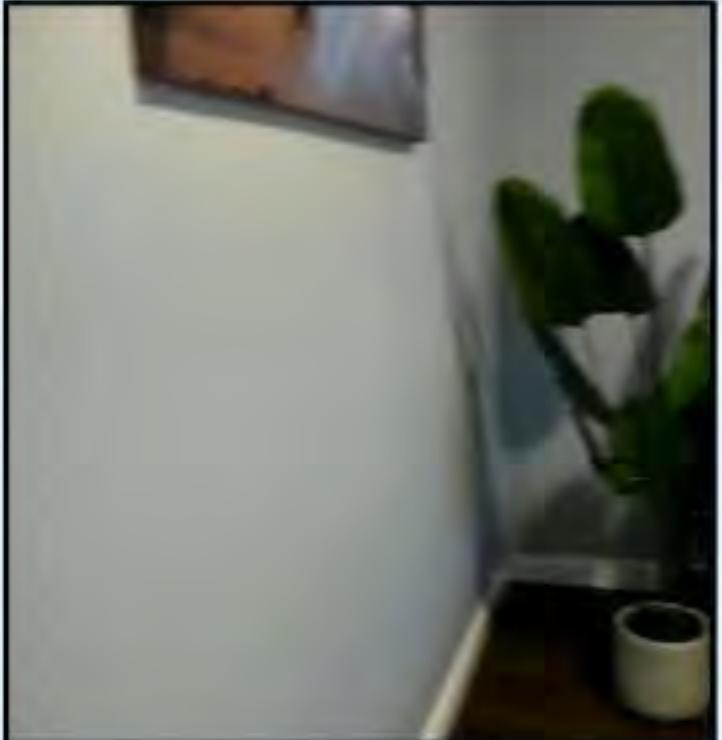


Baselines

Ours



Success



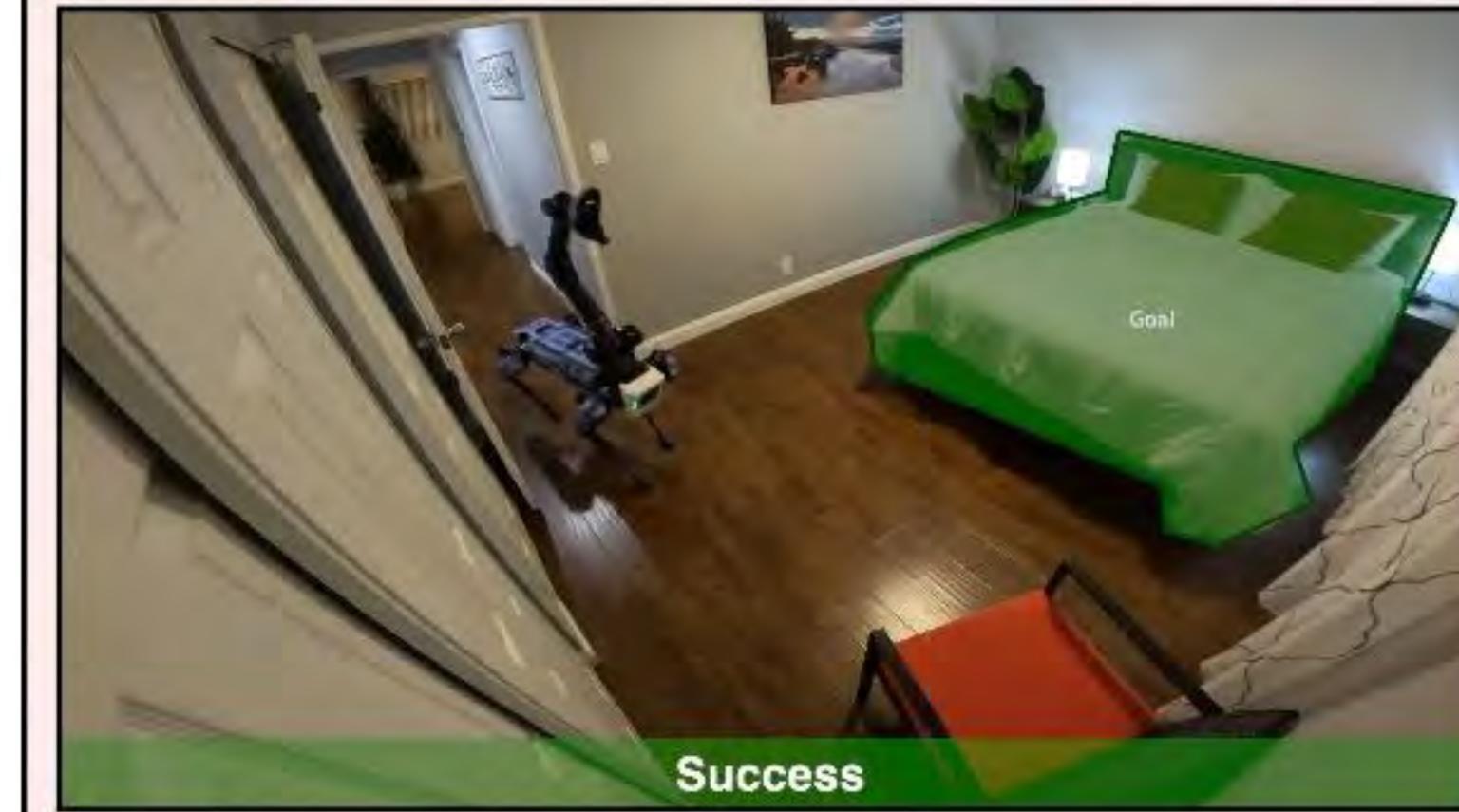
Observation



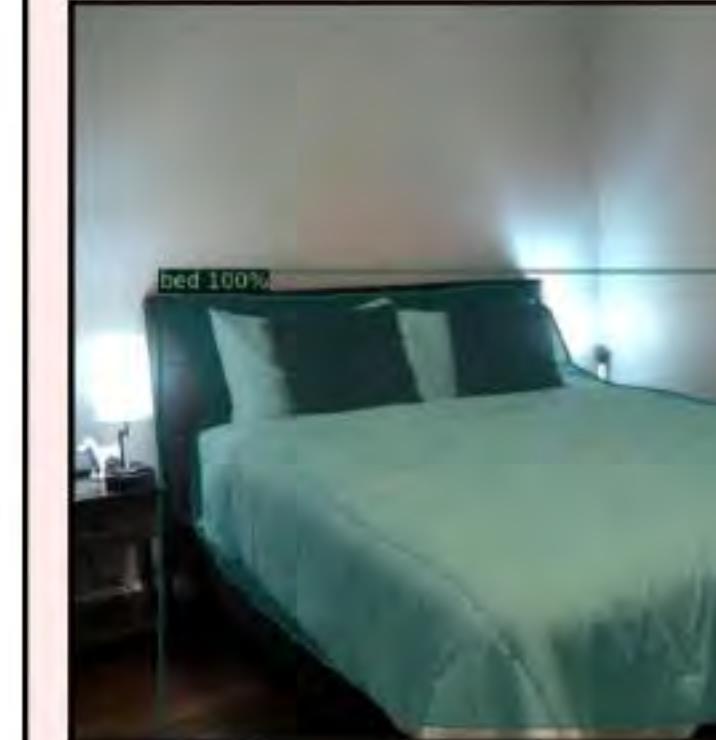
Instance Map

Success: 6/6 SPL: 0.78

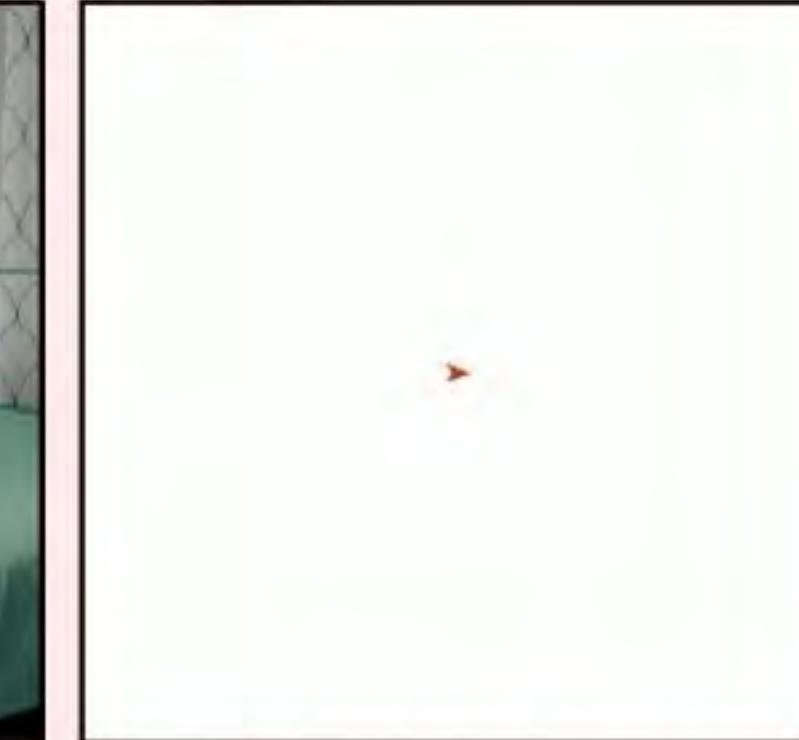
No Memory



Success



Observation



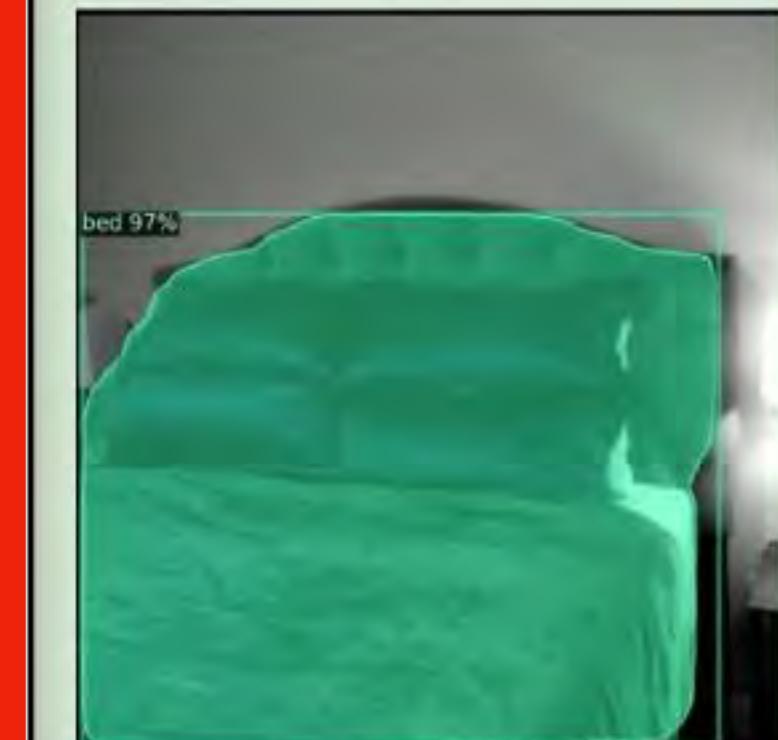
Instance Map

Success: 4/6 SPL: 0.40

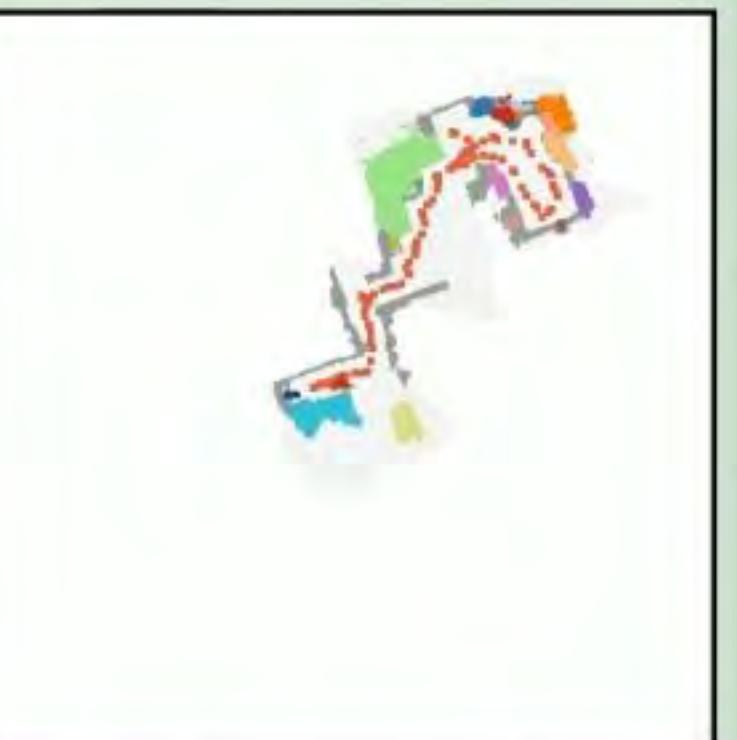
CoW



Failure



Observation



Instance Map

Success: 1/6 SPL: 0.16

Goal:

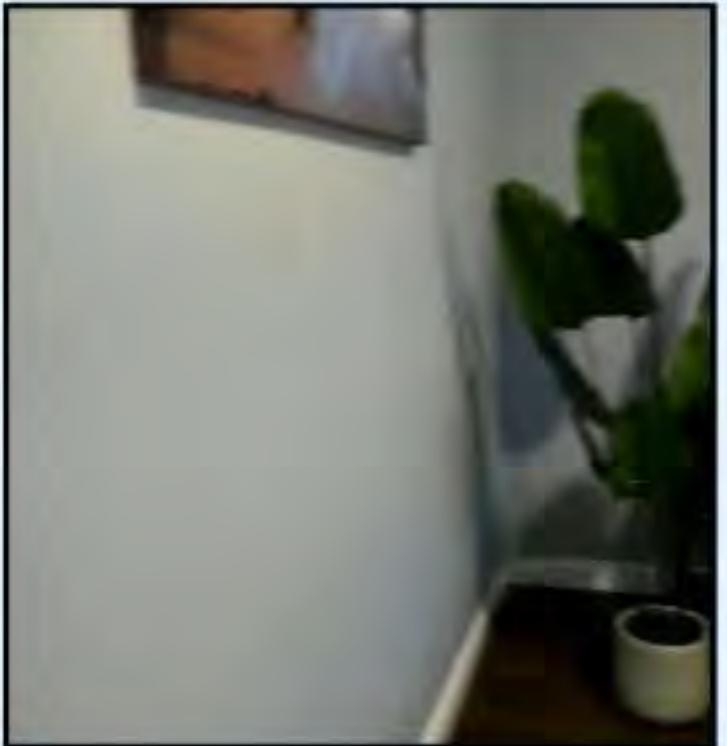


Baselines

Ours



Success



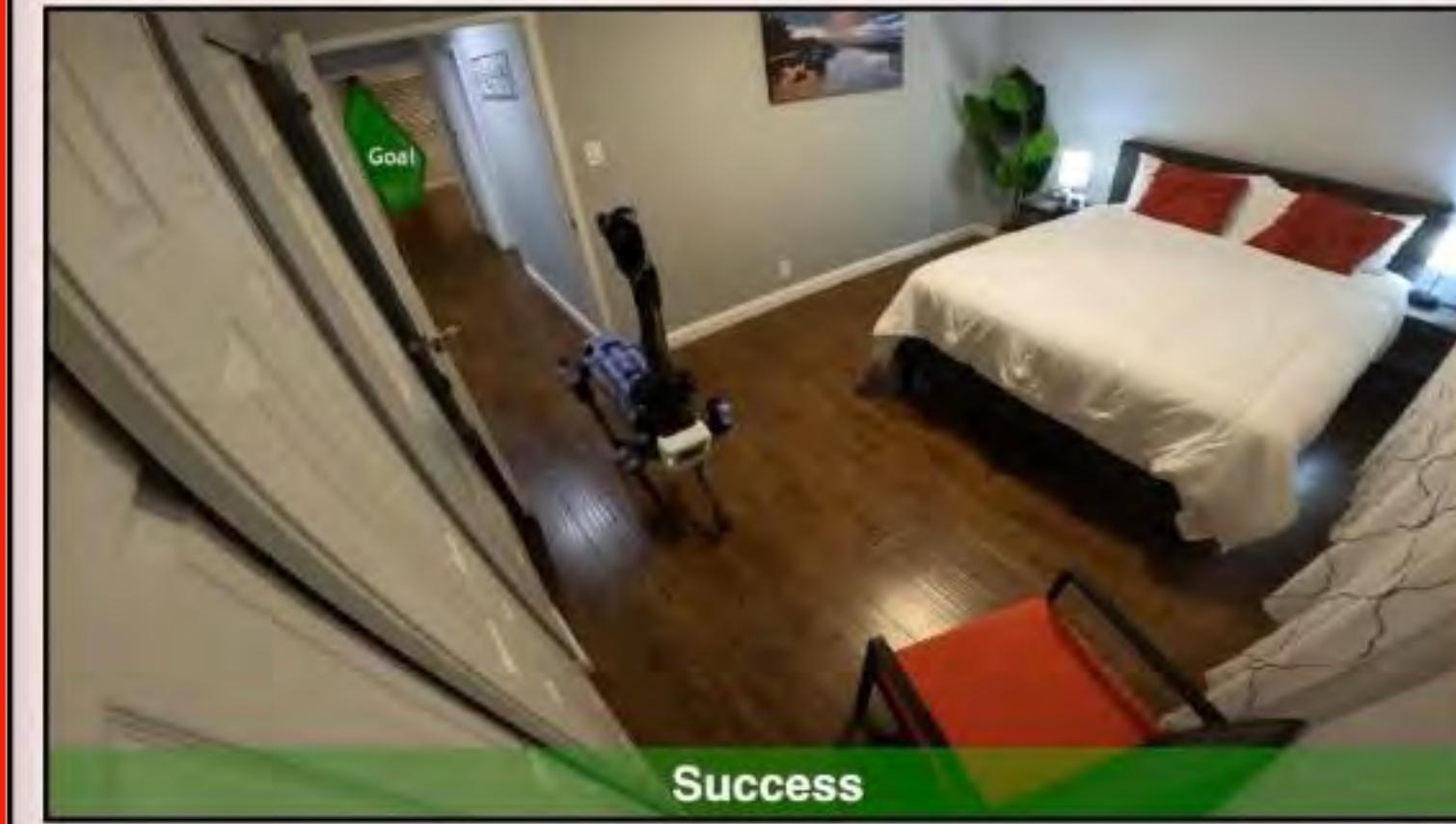
Observation



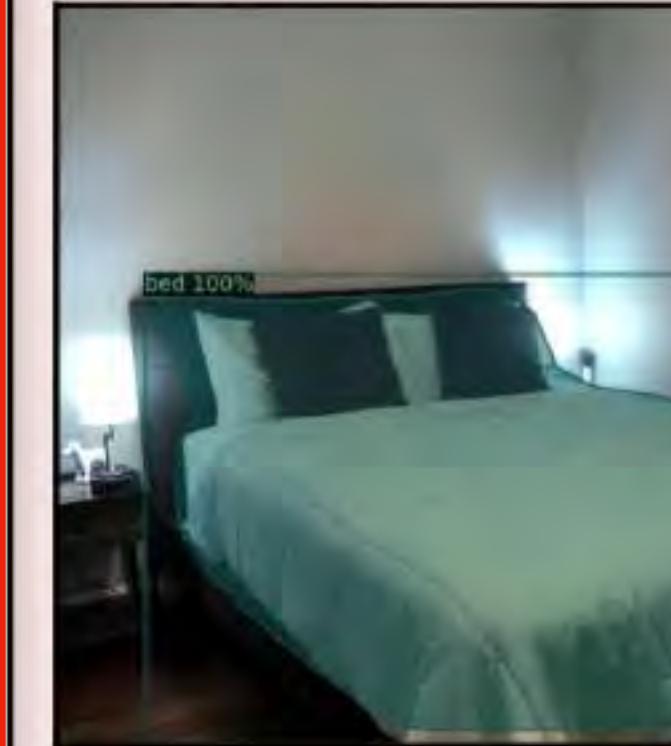
Instance Map

Success: 6/6 SPL: 0.78

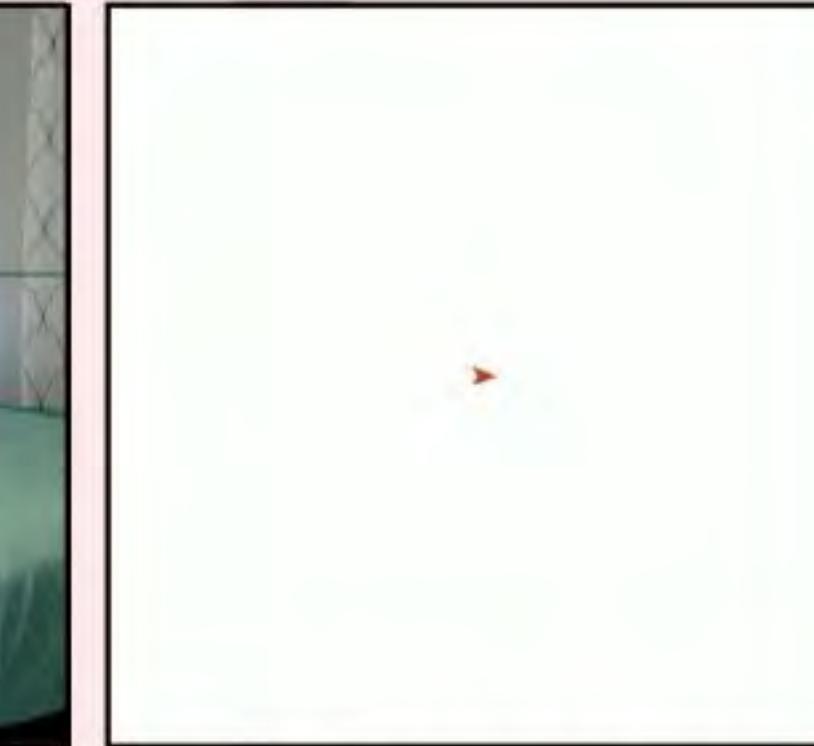
No Memory



Success



Observation



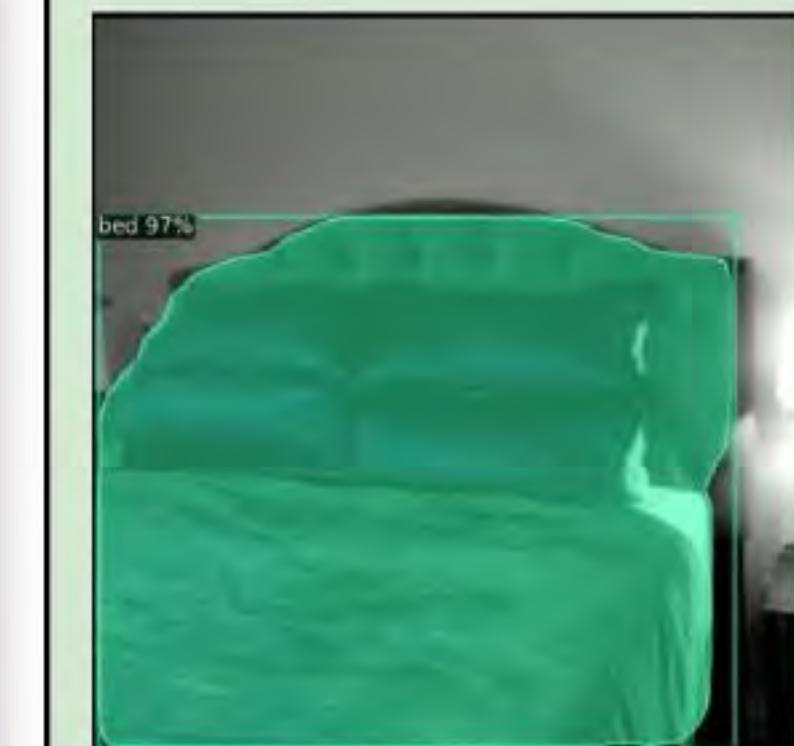
Instance Map

Success: 4/6 SPL: 0.40

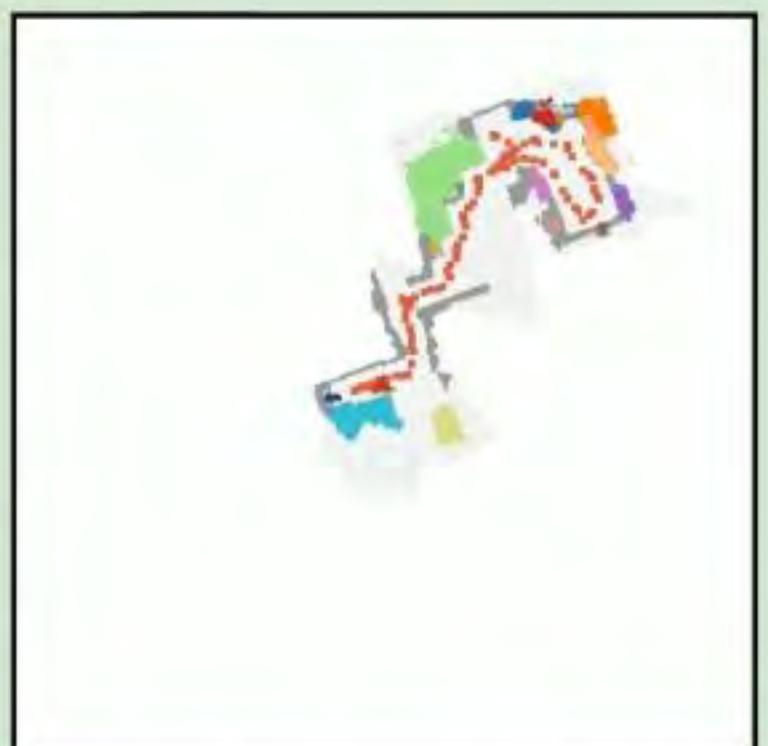
CoW



Failure



Observation



Instance Map

Success: 1/6 SPL: 0.16

Goal:

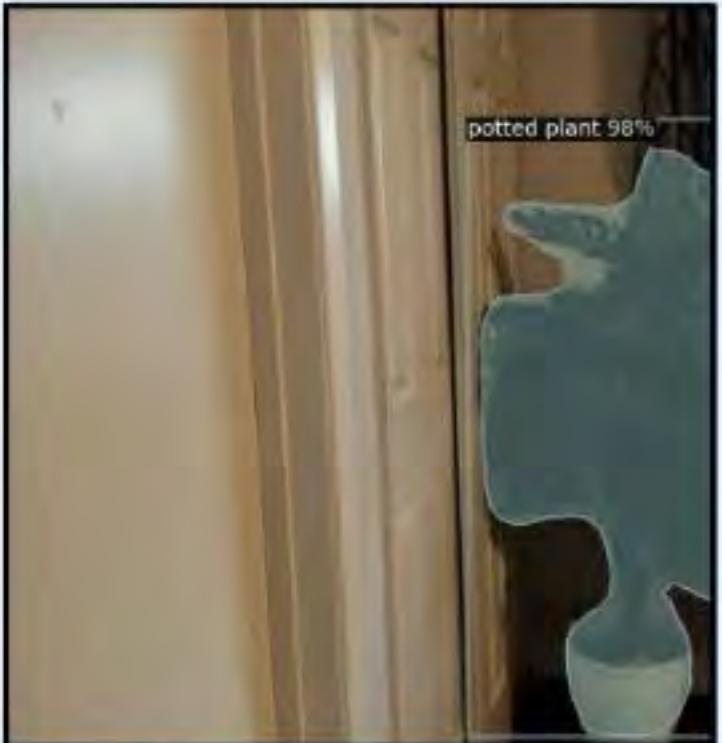


Baselines

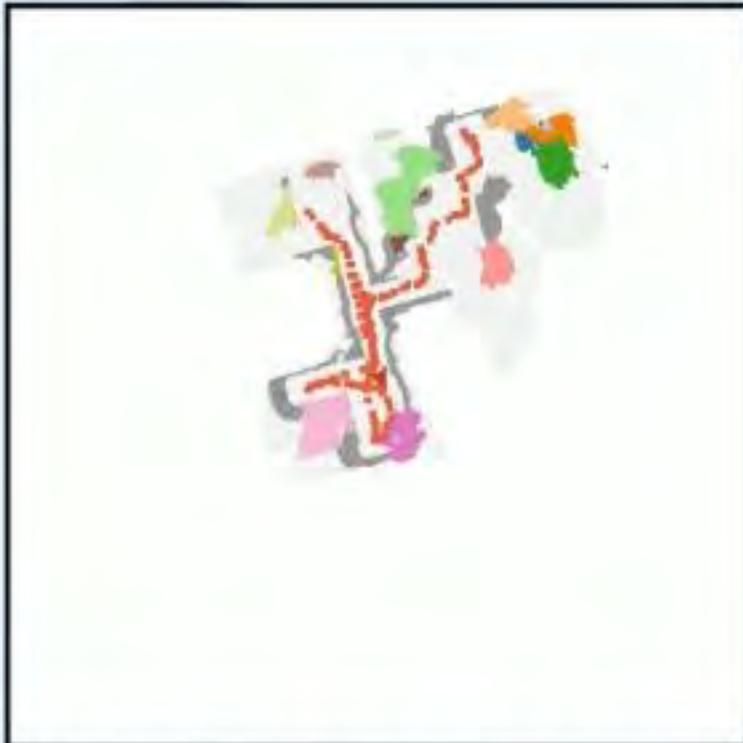
Ours



Success



Observation



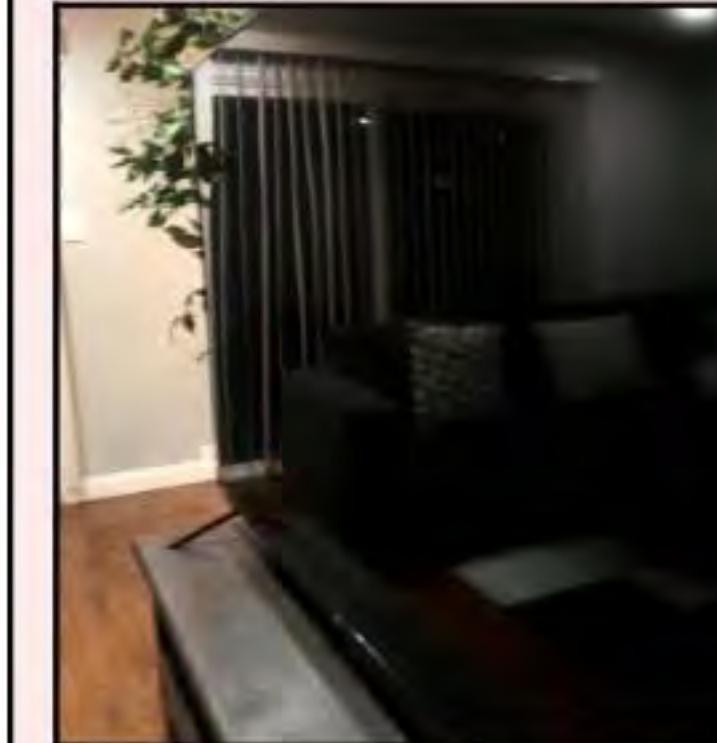
Instance Map

Success: 6/6 SPL: 0.78

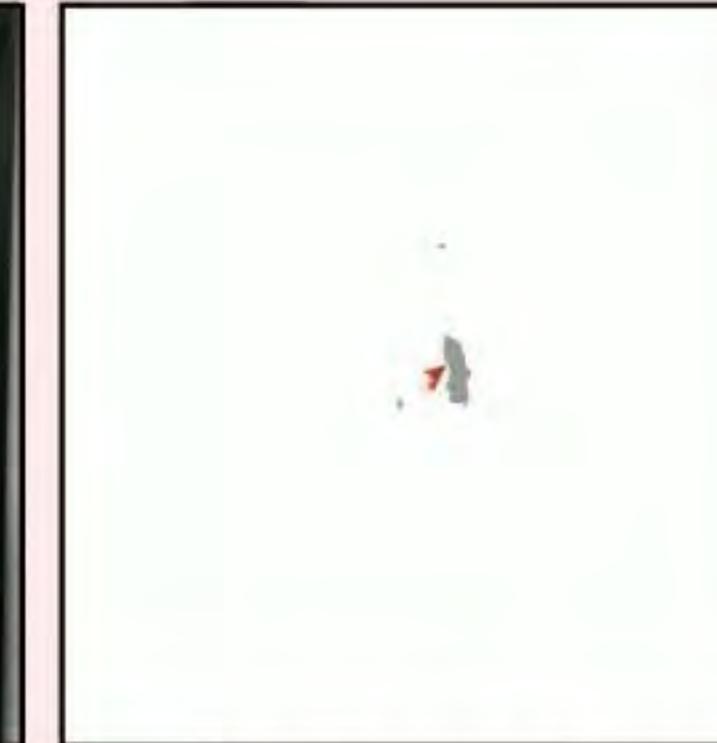
No Memory



Failure



Observation



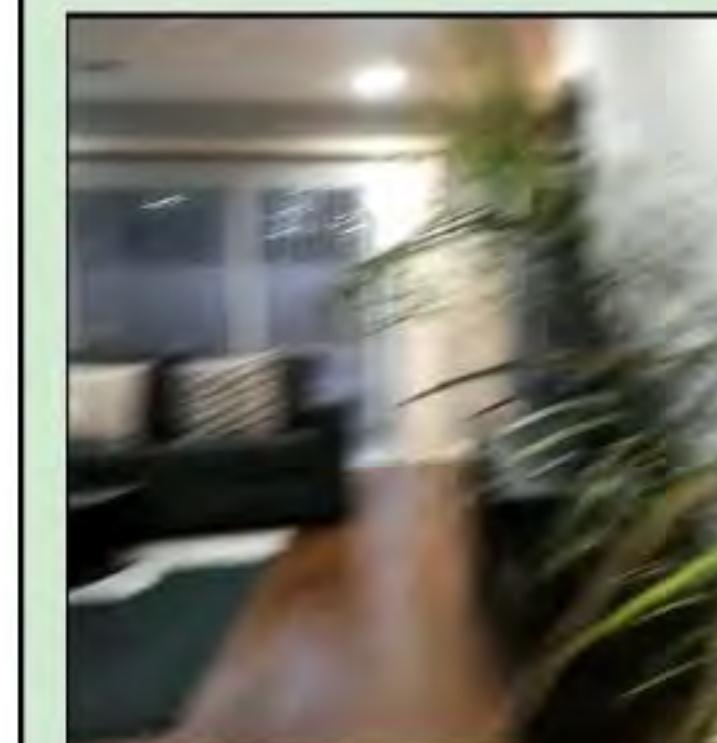
Instance Map

Success: 4/6 SPL: 0.40

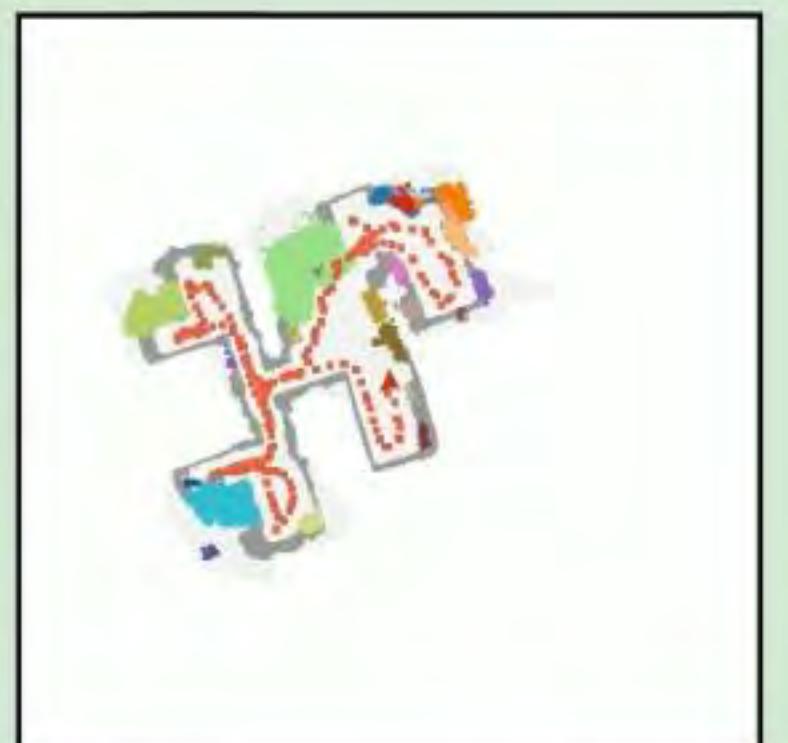
CoW



Failure



Observation



Instance Map

Success: 1/6 SPL: 0.16

1. GOAT Problem

GOAT System Architecture

Results

Applications

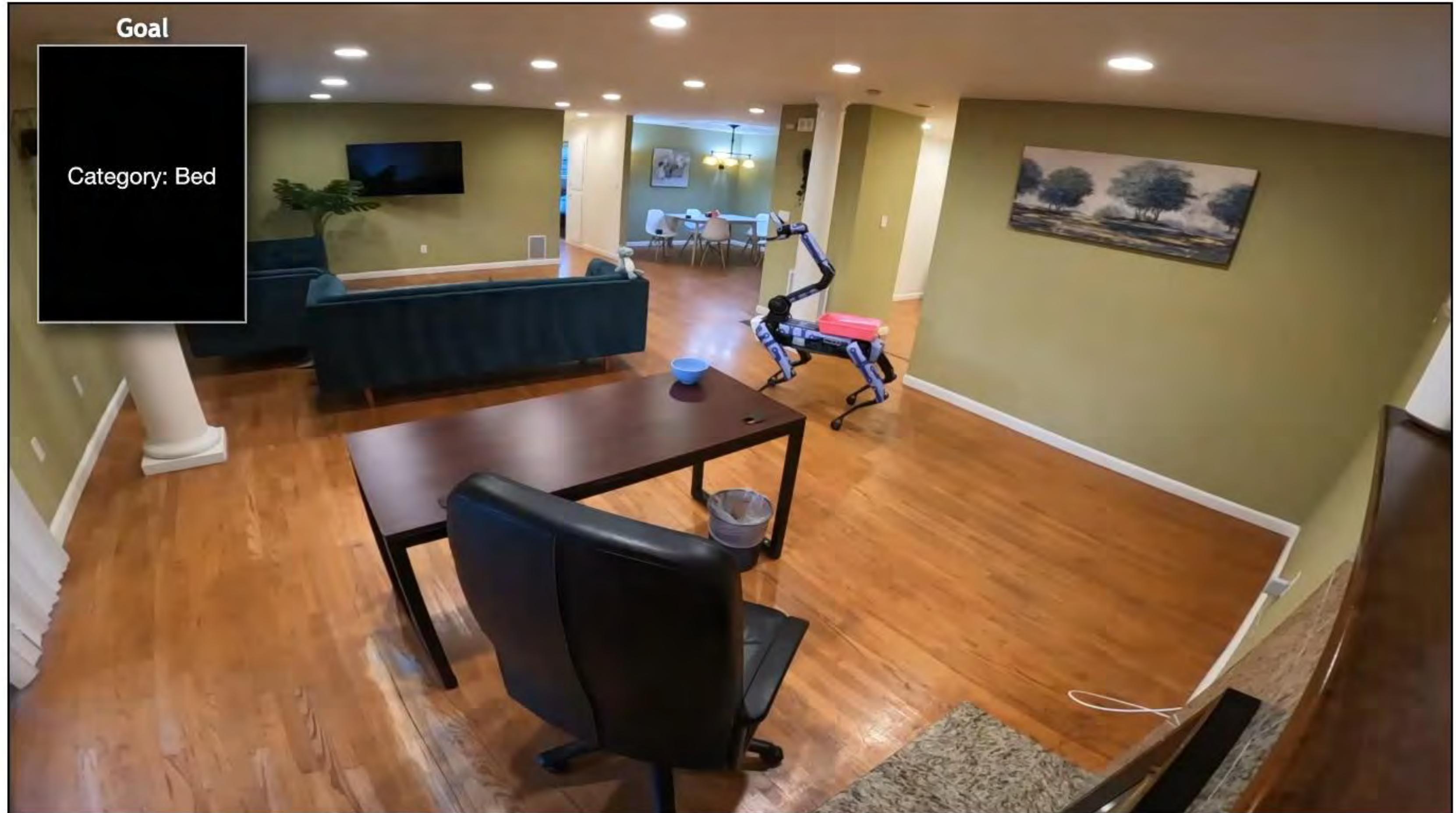
Pick & Place

Social Navigation

Platform Agnostic

Pick & Place

Third-person view



Social Navigation

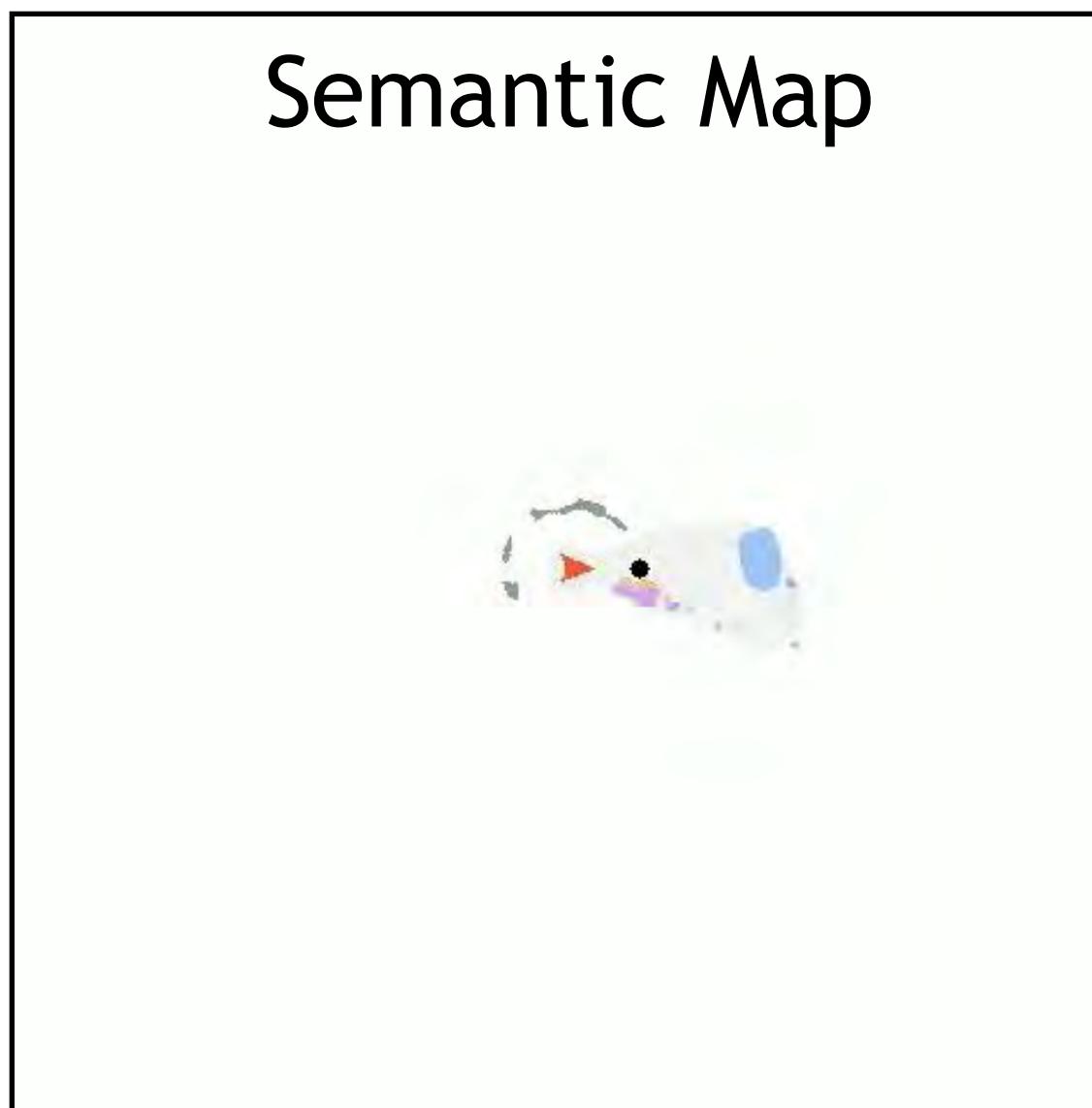
Observation



Third-person view



Semantic Map



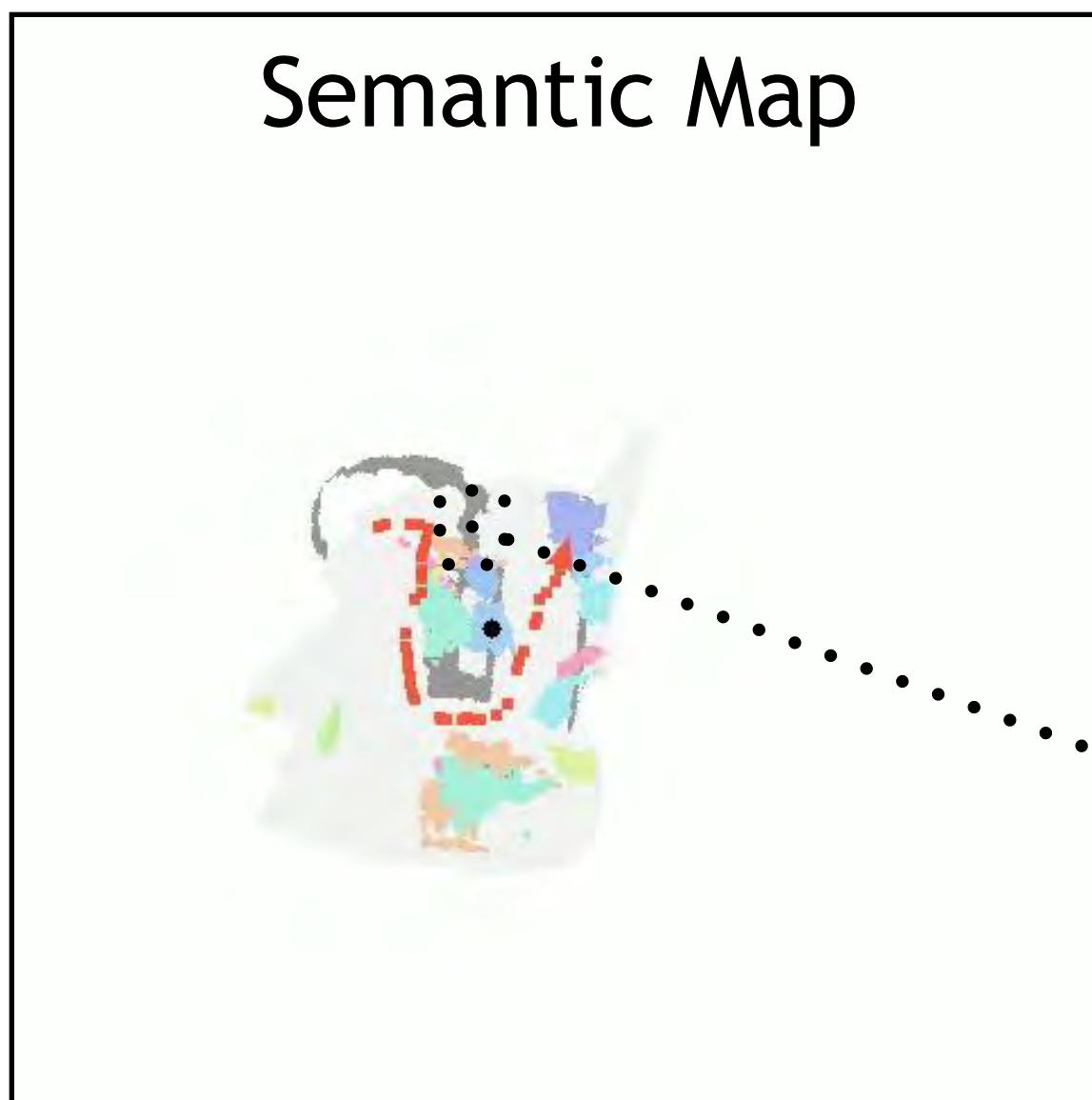
Robot plans around the dynamic obstacle (person) to go to the refrigerator

Social Navigation

Observation



Semantic Map



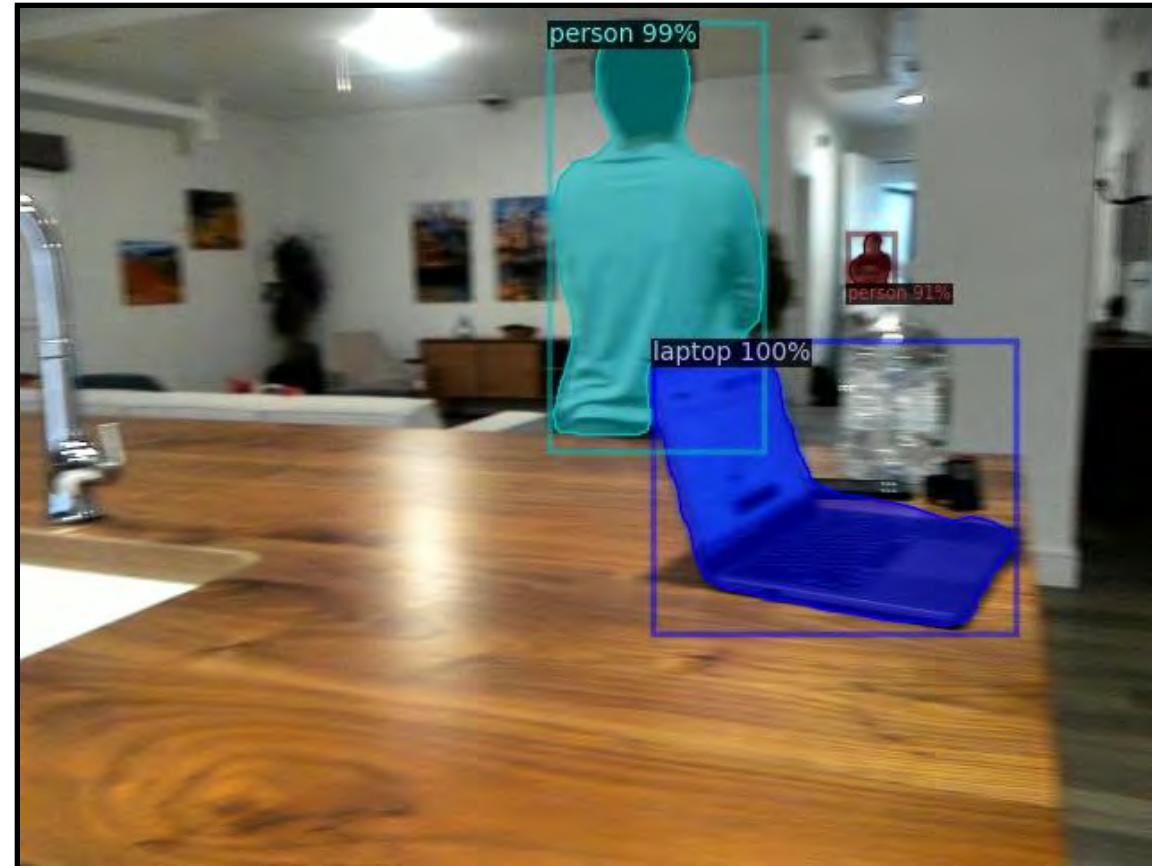
Third-person view



Robot removes previous location of person from the map

Social Navigation

Observation



Semantic Map



Third-person view



Robot follows the person while updating their location

Platform Agnostic

Third-person view



Summary

Universal navigation

- Multimodal

Image



Language

*Find the fruit
basket on the
kitchen
counter*

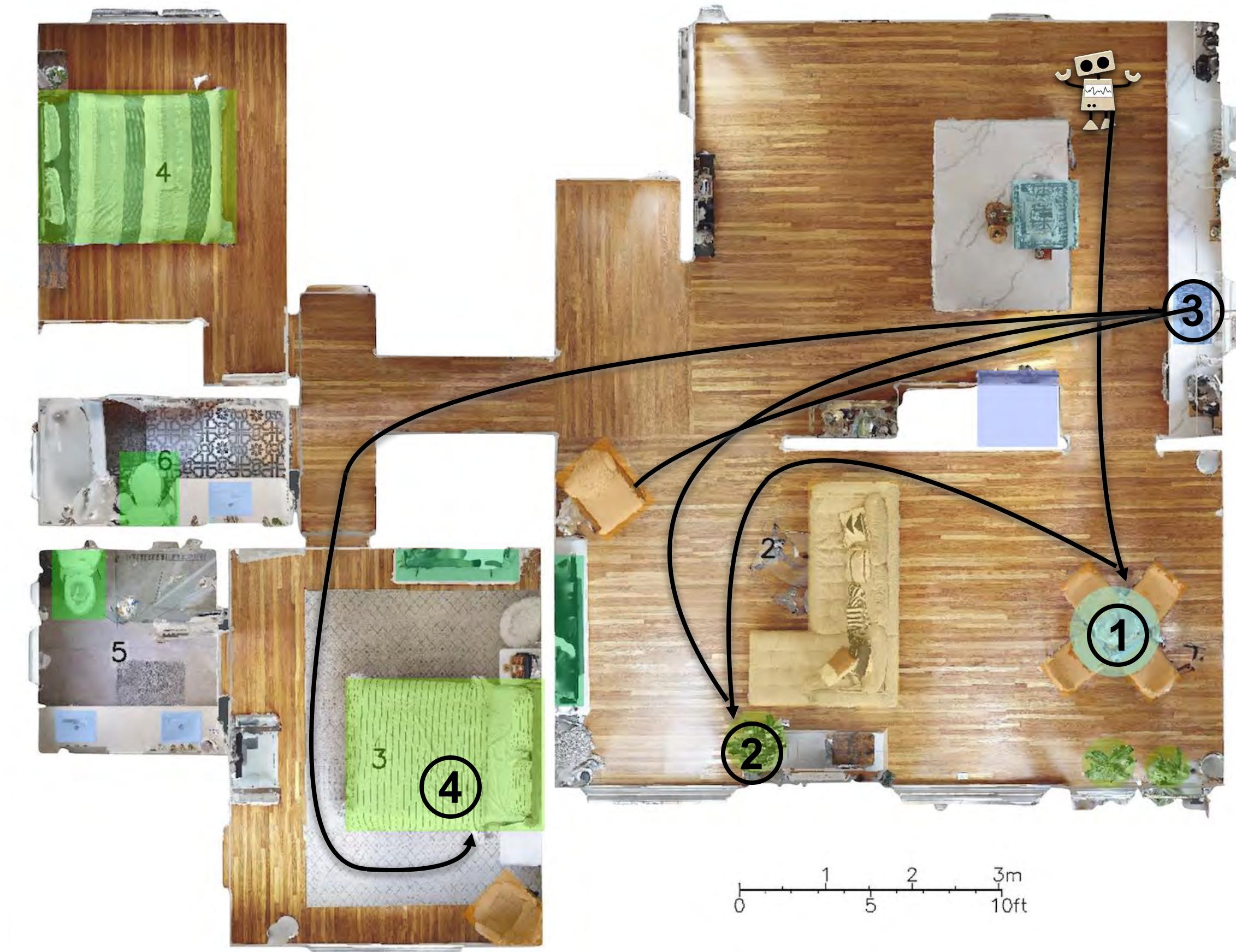
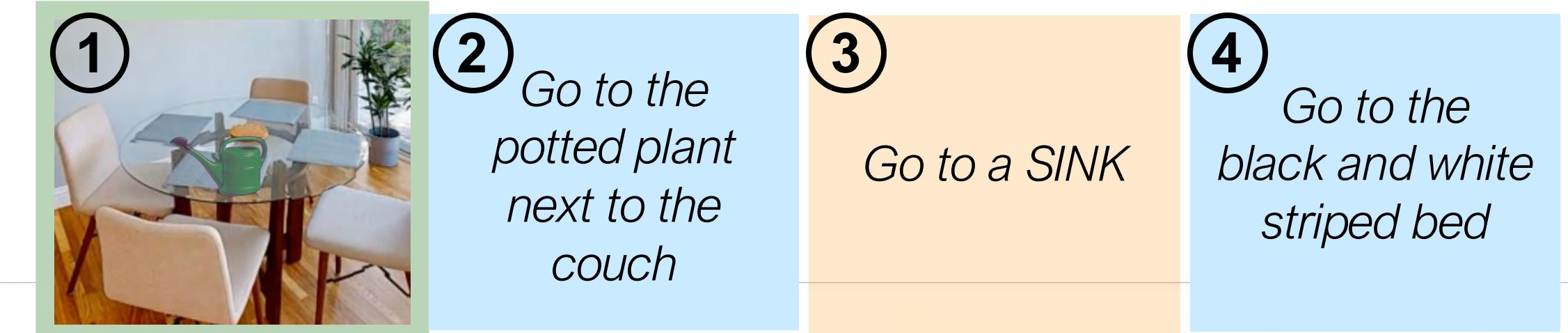
Category

*Bring me
a CUP*

Summary

Universal navigation

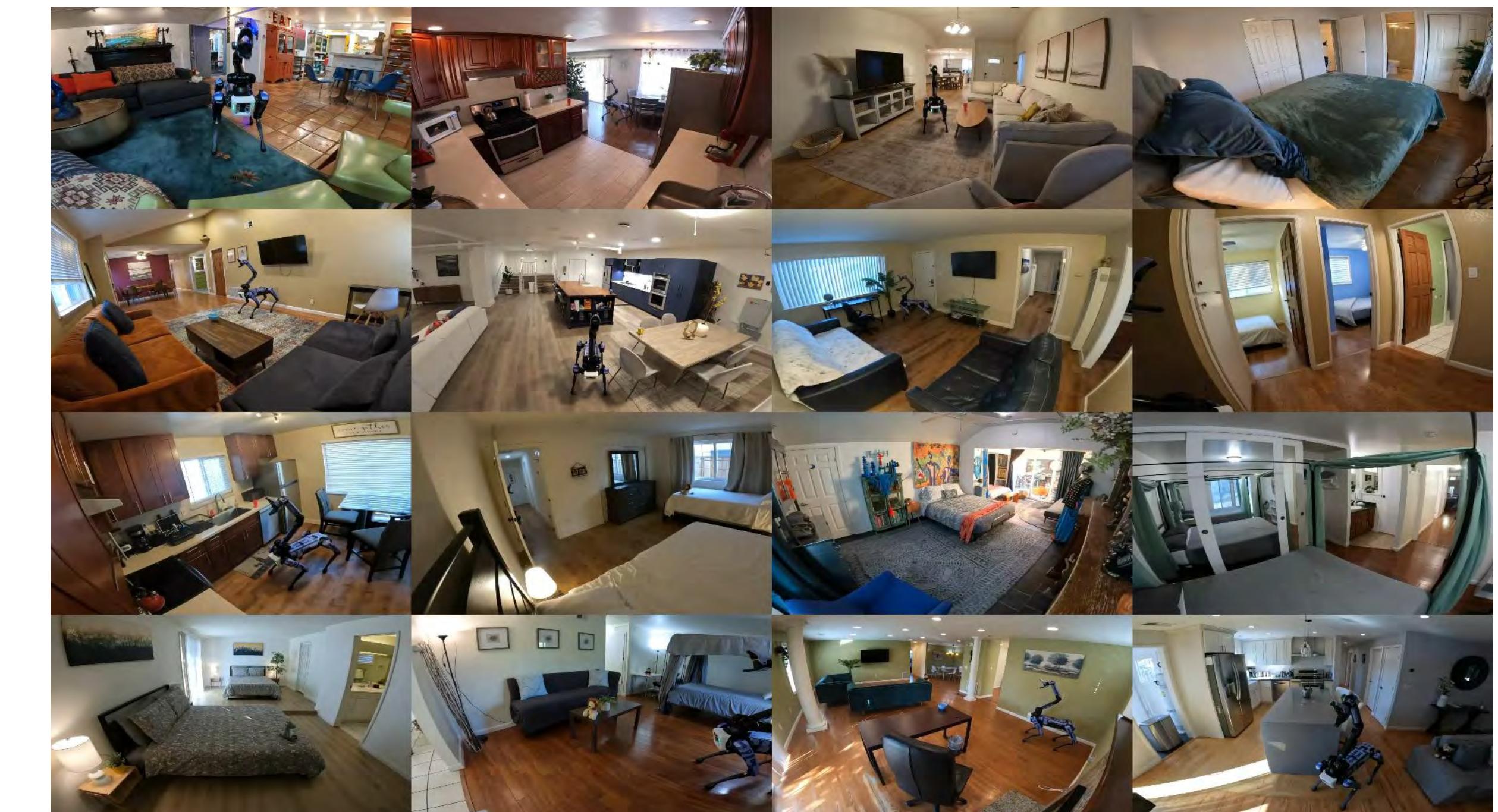
- Multimodal
- Lifelong



Summary

Universal navigation

- Multimodal
- Lifelong
- Unseen environments



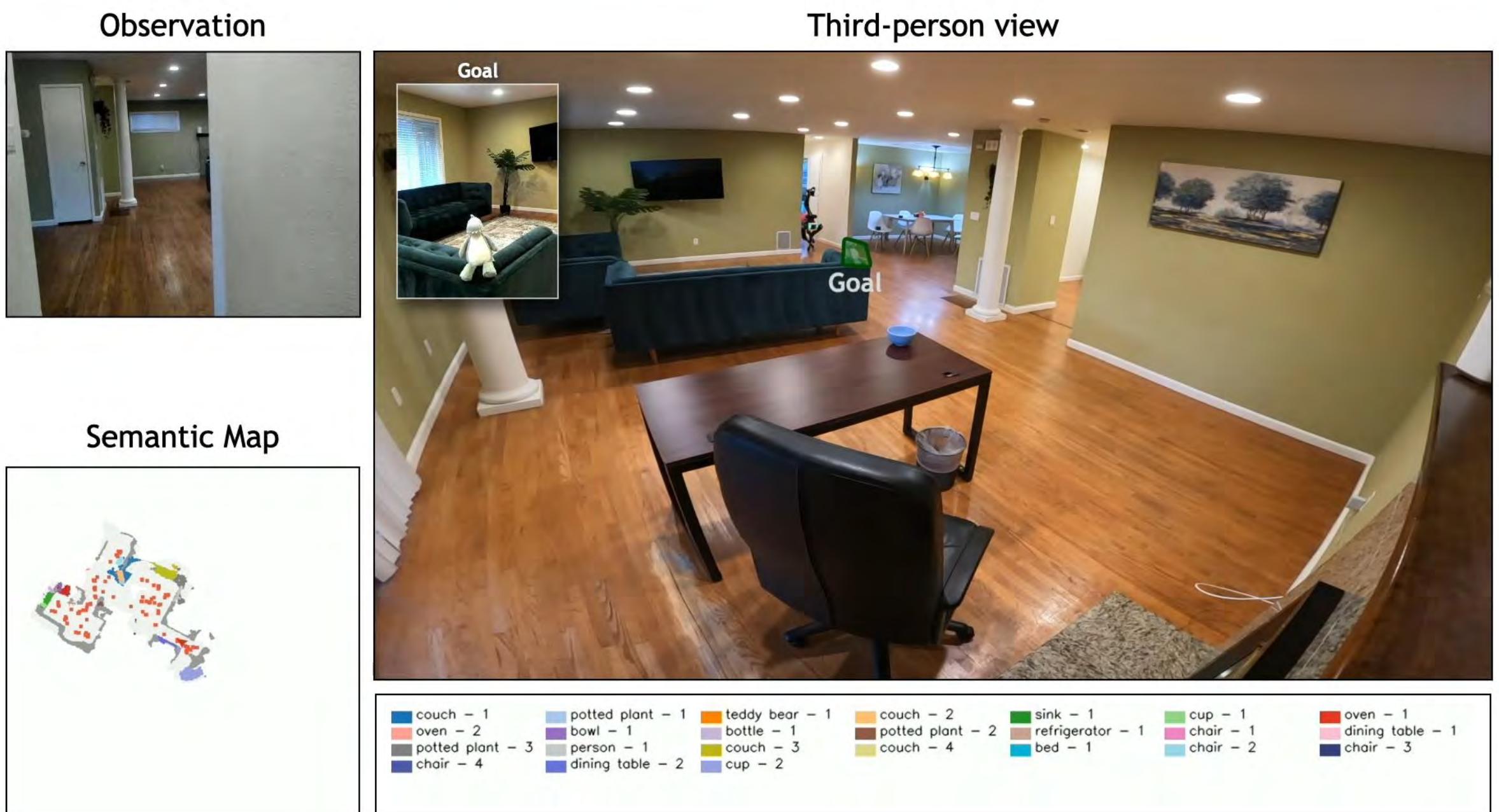
Summary

Universal navigation

- Multimodal
- Lifelong
- Unseen environments

Applications

- Pick & Place



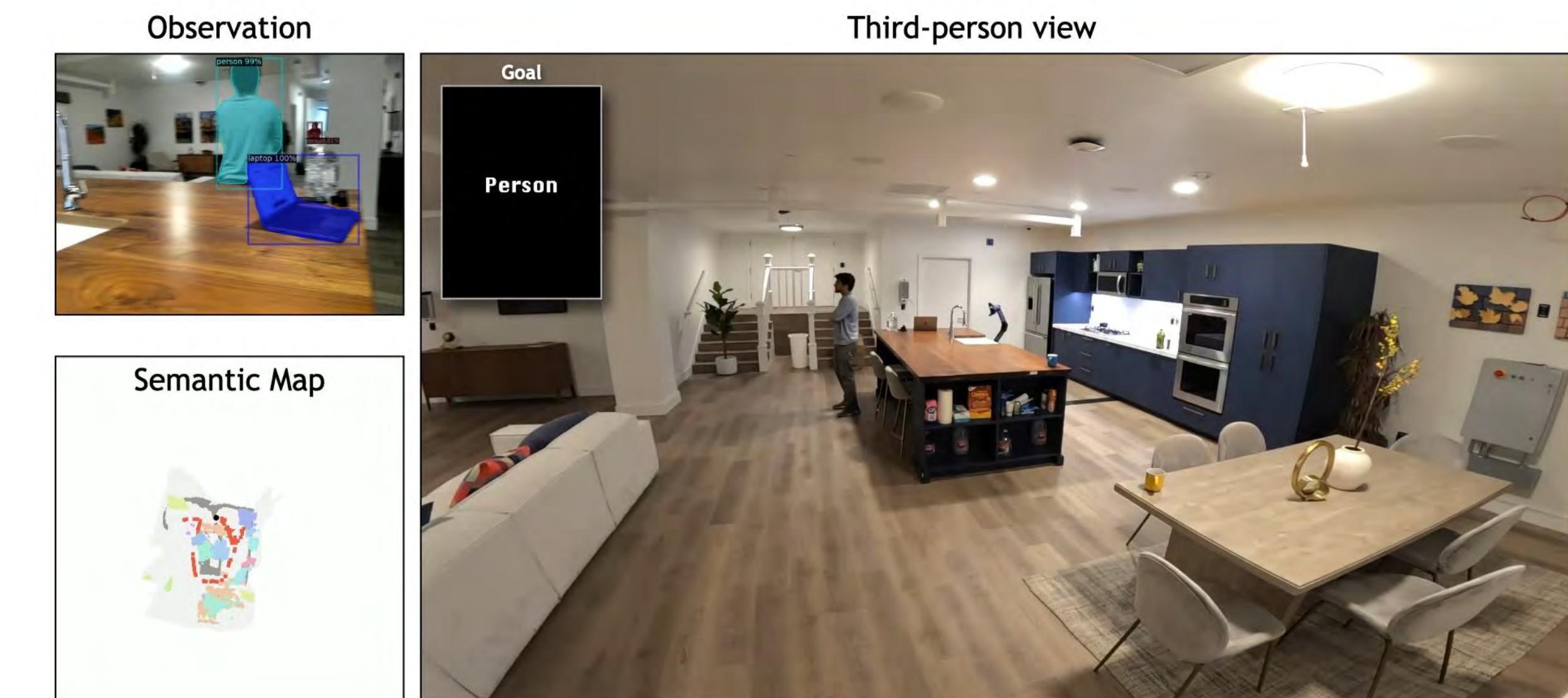
Summary

Universal navigation

- Multimodal
- Lifelong
- Unseen environments

Applications

- Pick & Place
- Social Navigation



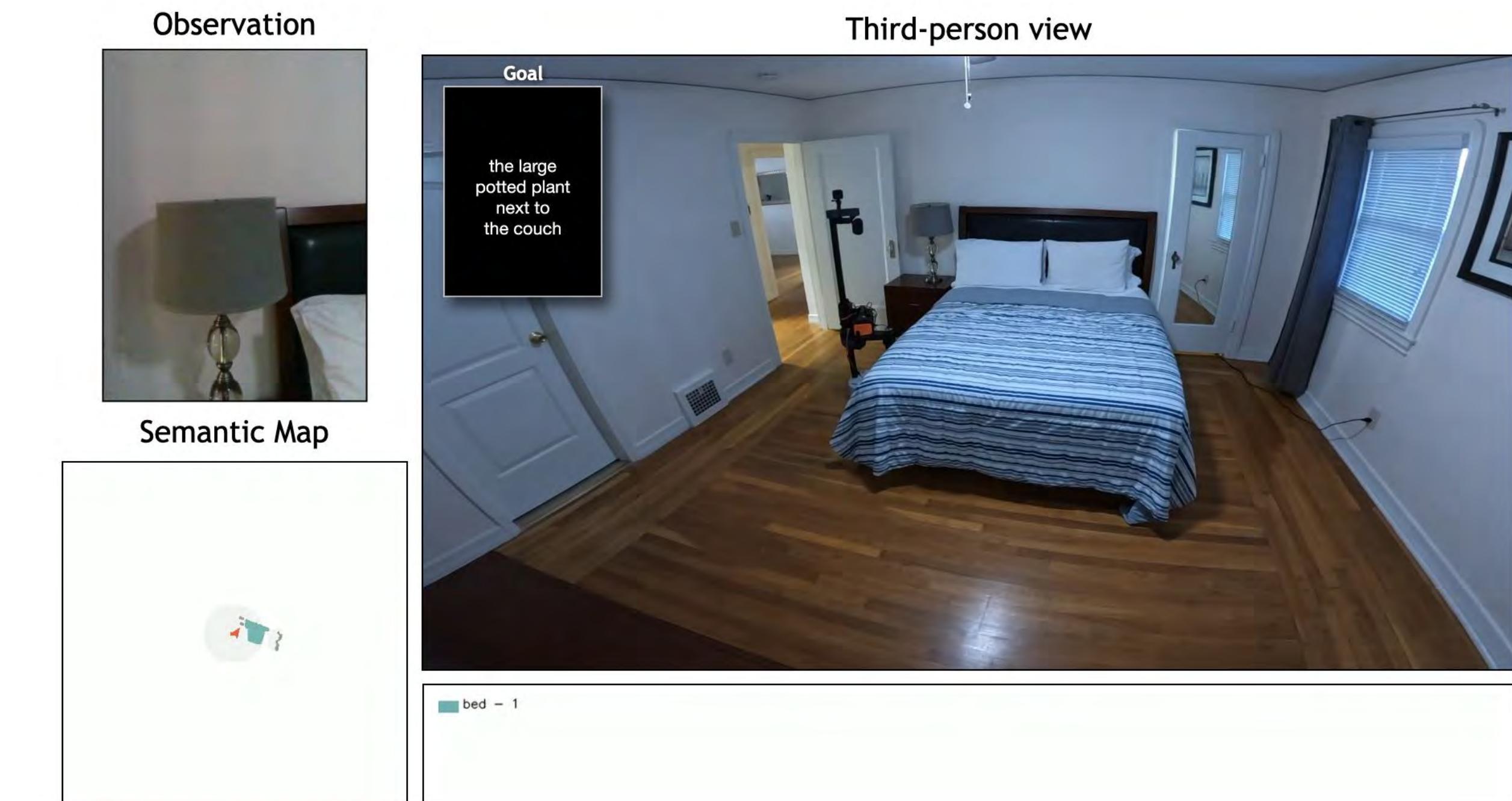
Summary

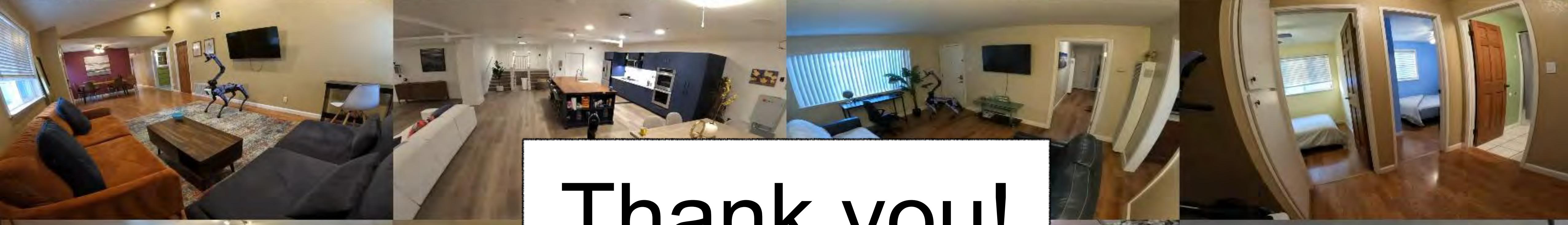
Universal navigation

- Multimodal
- Lifelong
- Unseen environments

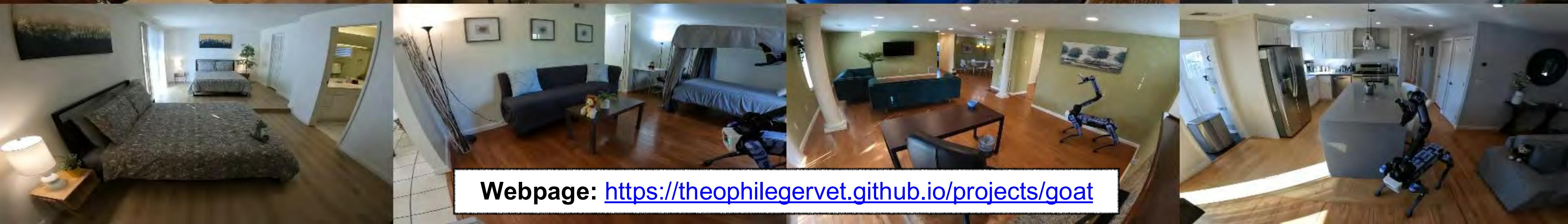
Applications

- Pick & Place
- Social Navigation
- Platform Agnostic





Thank you!



Webpage: <https://theophilegervet.github.io/projects/goat>

Thank you!

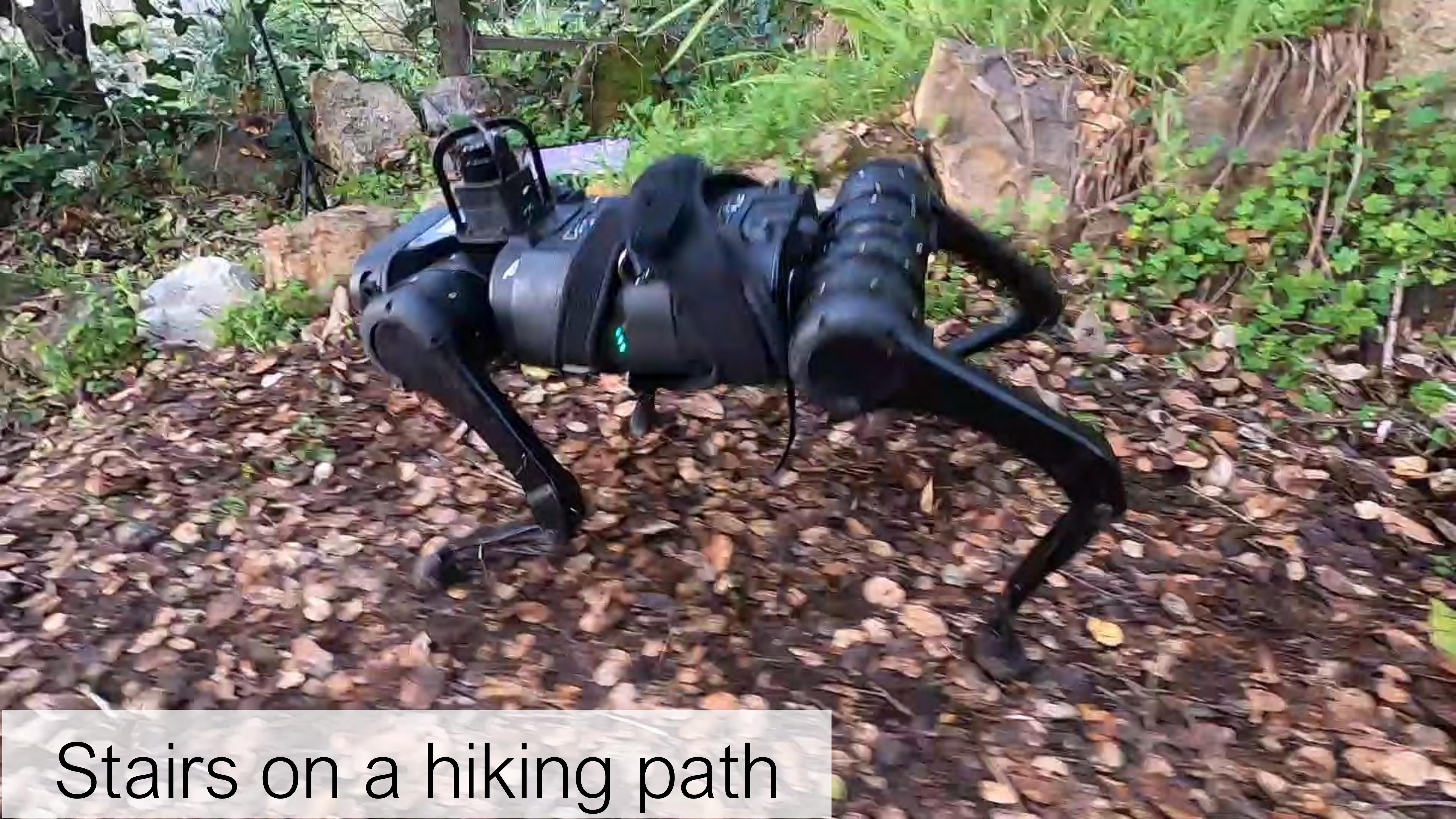
Learning to Walk for Quadrupeds

Jitendra Malik
UC Berkeley



Recovery from leg obstruction

Rocky area next to river bed



Stairs on a hiking path

Unstable and constantly deforming ground



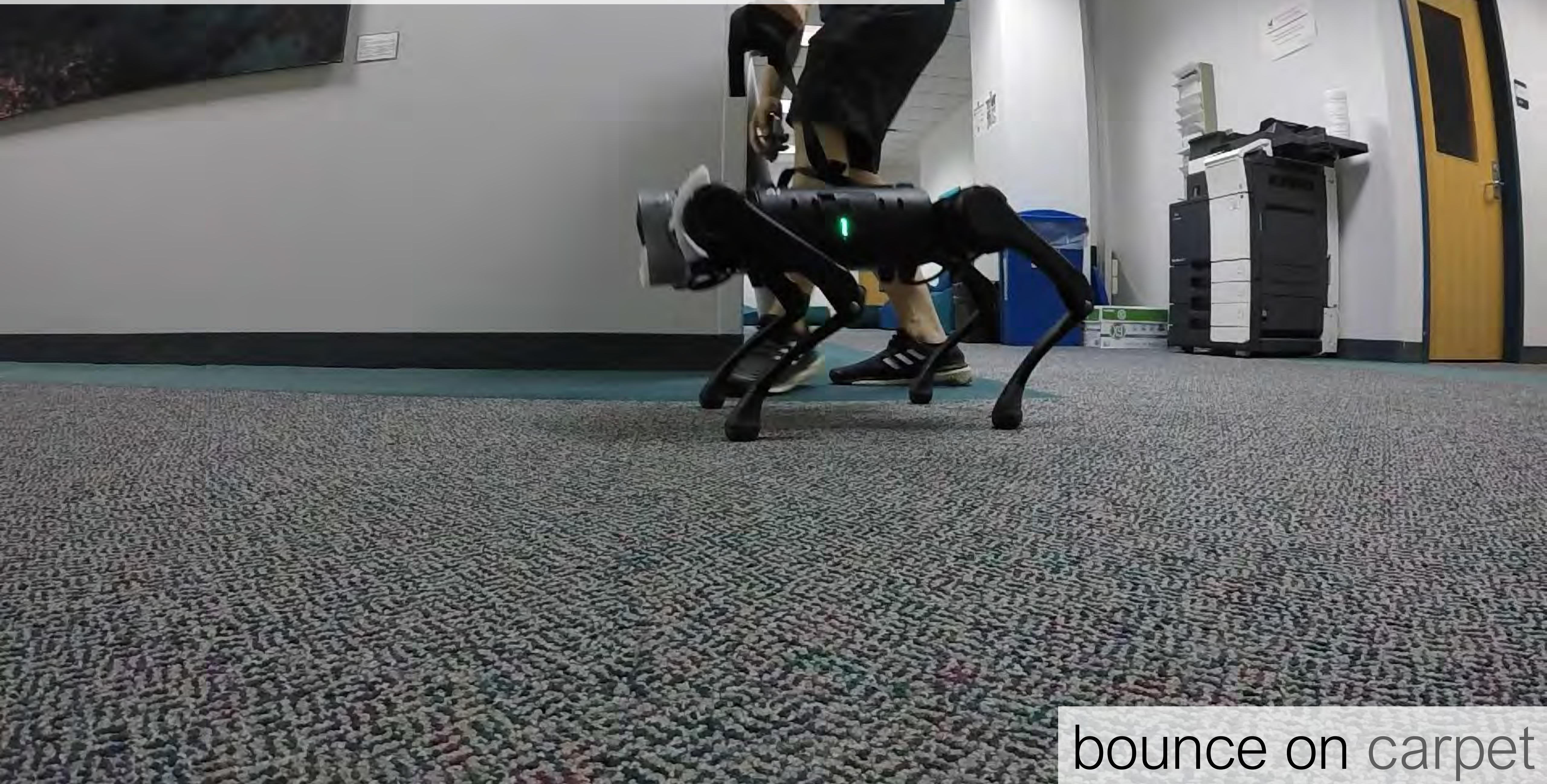
Loose Mud Pile at Construction Site



Vegetation on uneven surface



Bouncing (Gallop) Gait @ 1.5 m/s



bounce on carpet

Navigation in Cluttered Indoor Setting

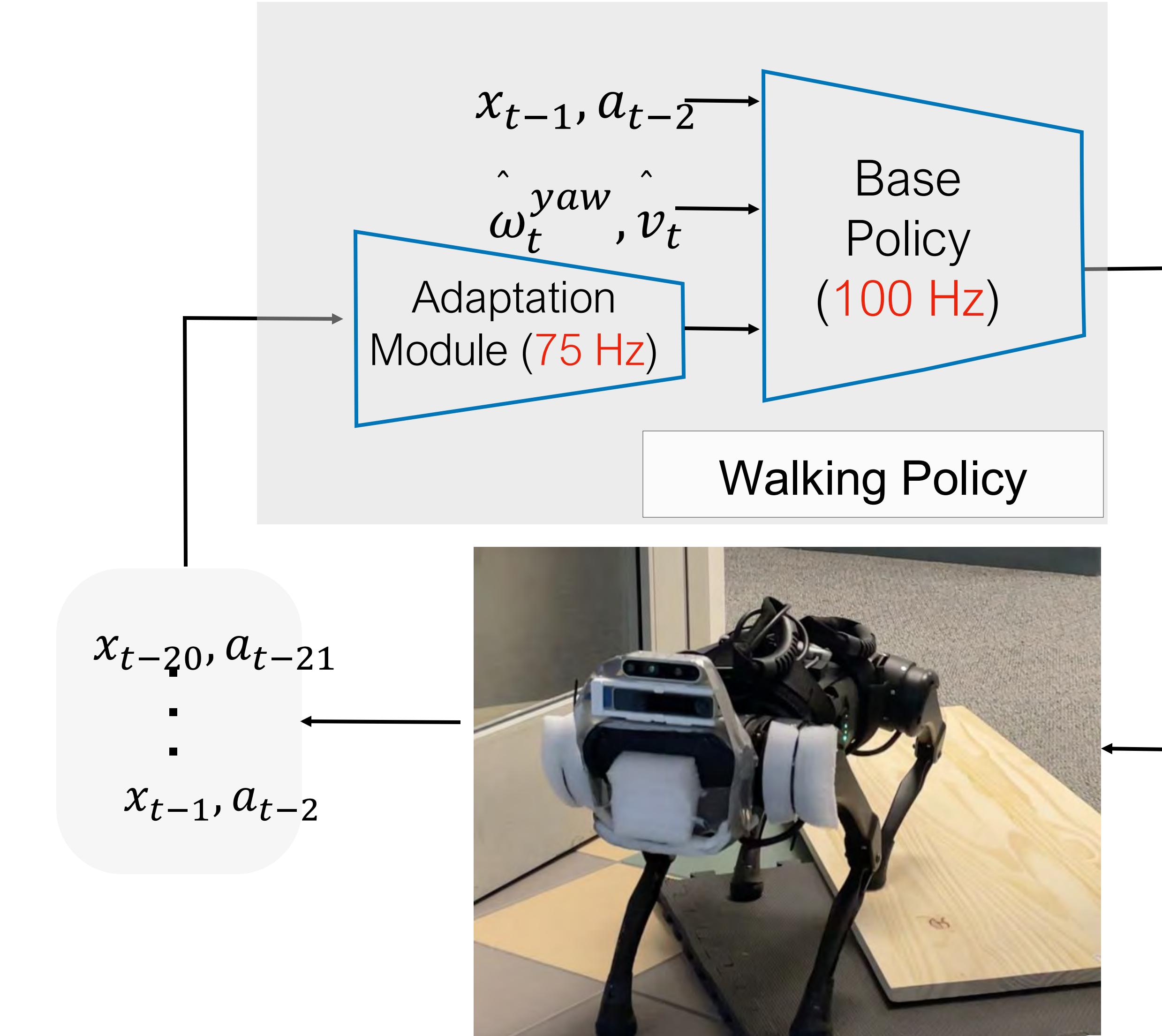


Coupling Vision and Proprioception for Navigation of Legged Robots

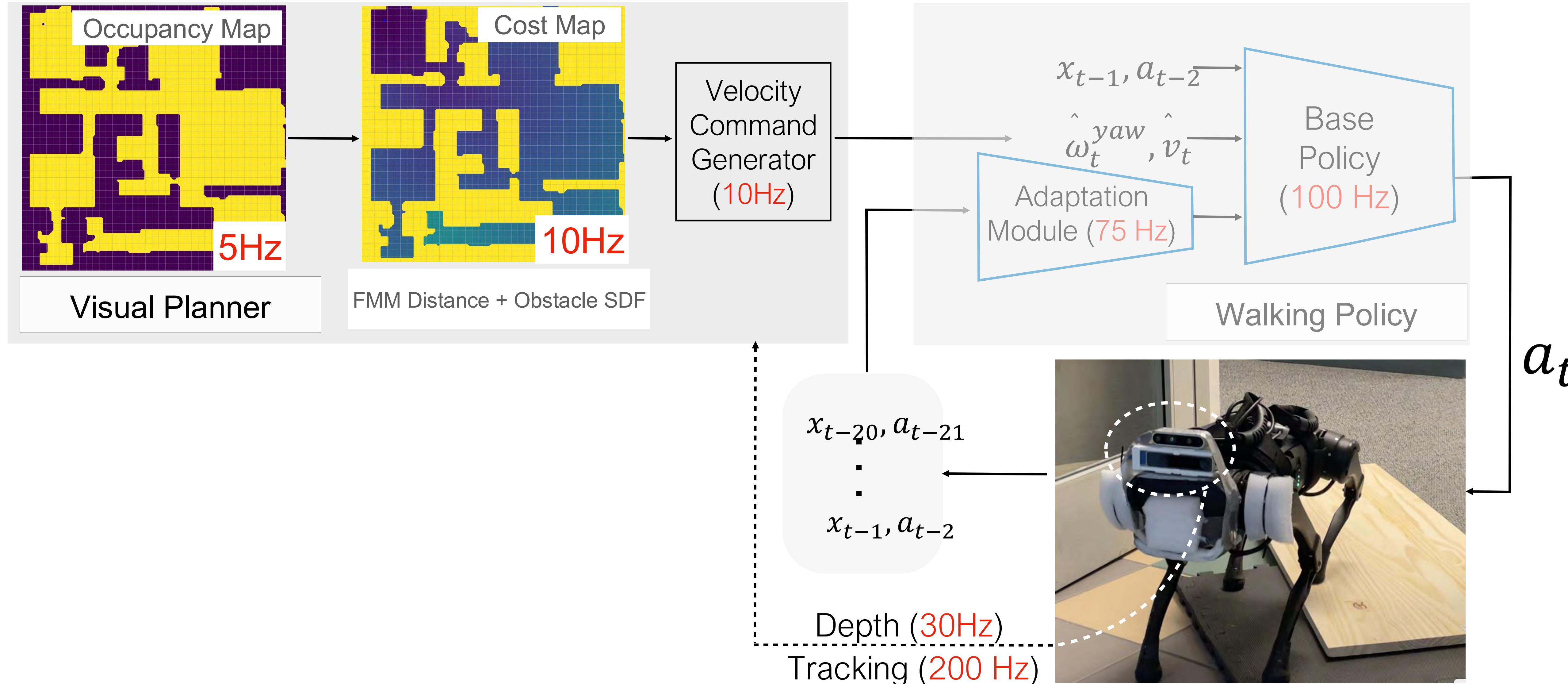
¹Zipeng Fu*, ²Ashish Kumar*, ¹Ananye Agarwal, ²Haozhi Qi,
²Jitendra Malik, ¹Deepak Pathak

¹CMU ²UC Berkeley
CVPR 2022

Linear and Angular Velocity



Visual Path Planner



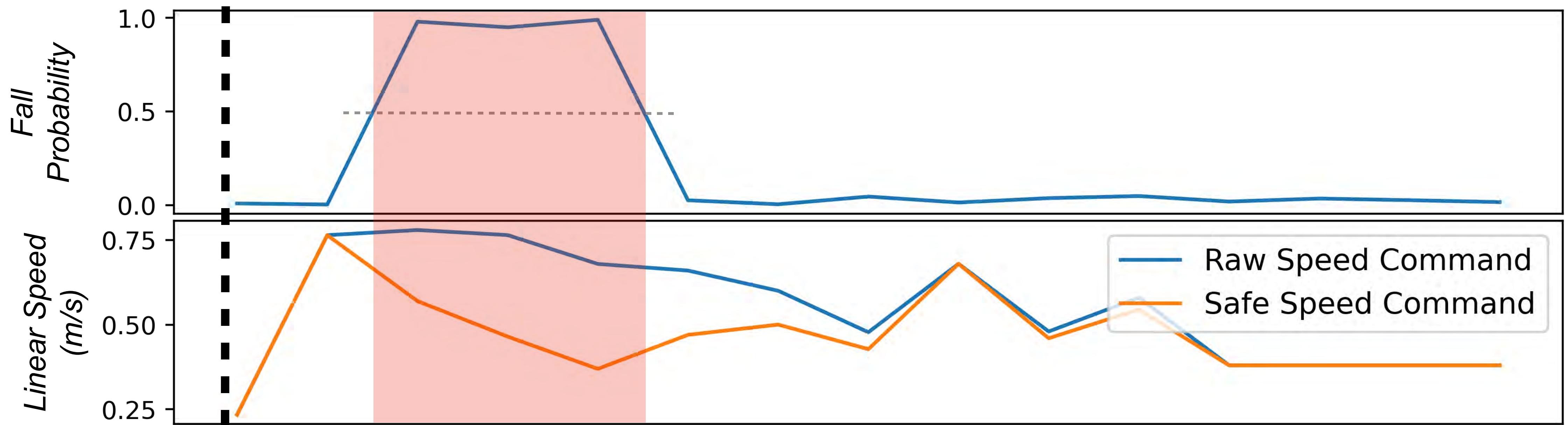
Navigation in Cluttered Indoor Setting



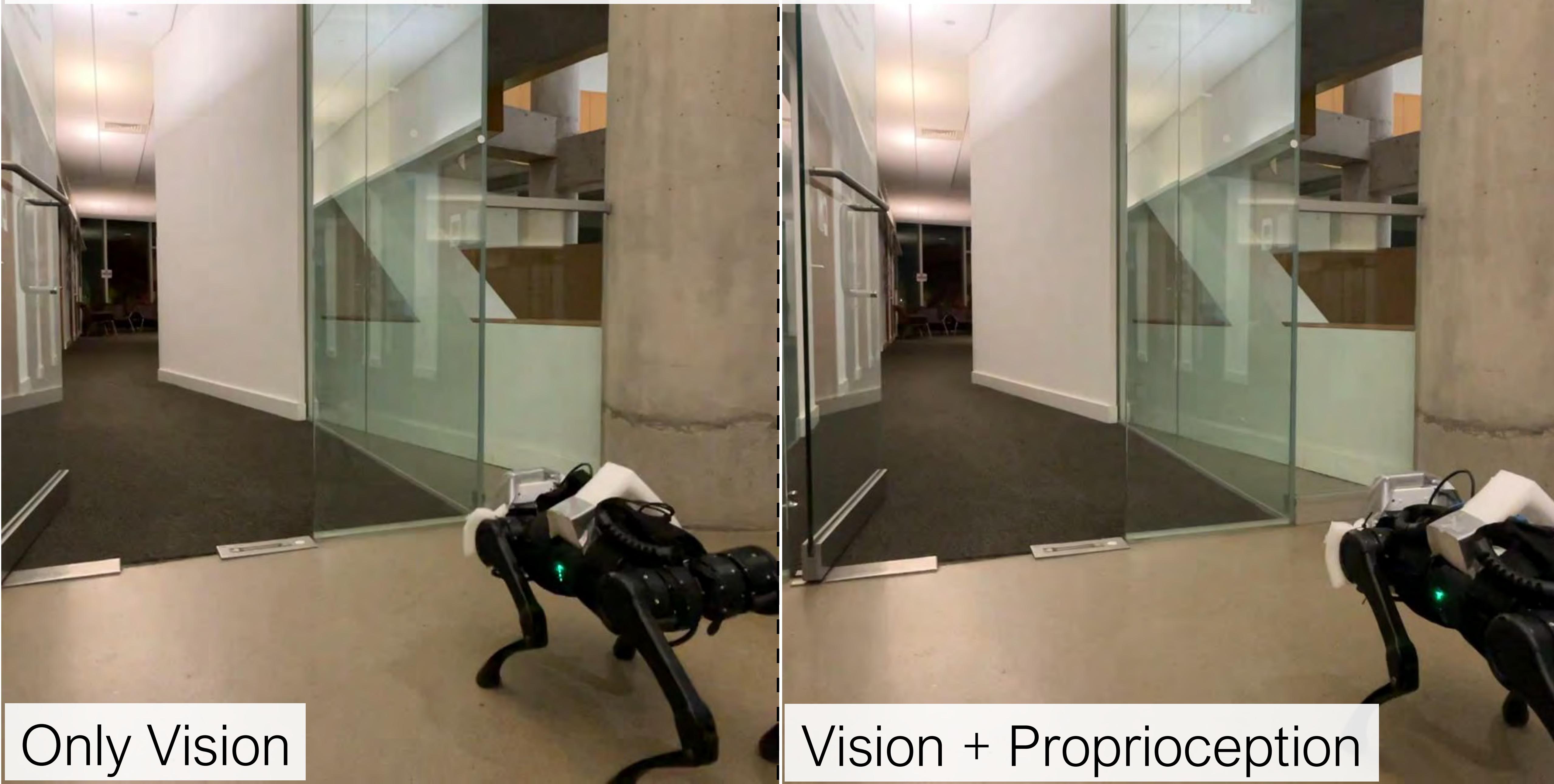
Whats Missing?

1. Navigation should be coordinated with locomotion
2. Vision should be coordinated with proprioception

Navigation coordinated with locomotion



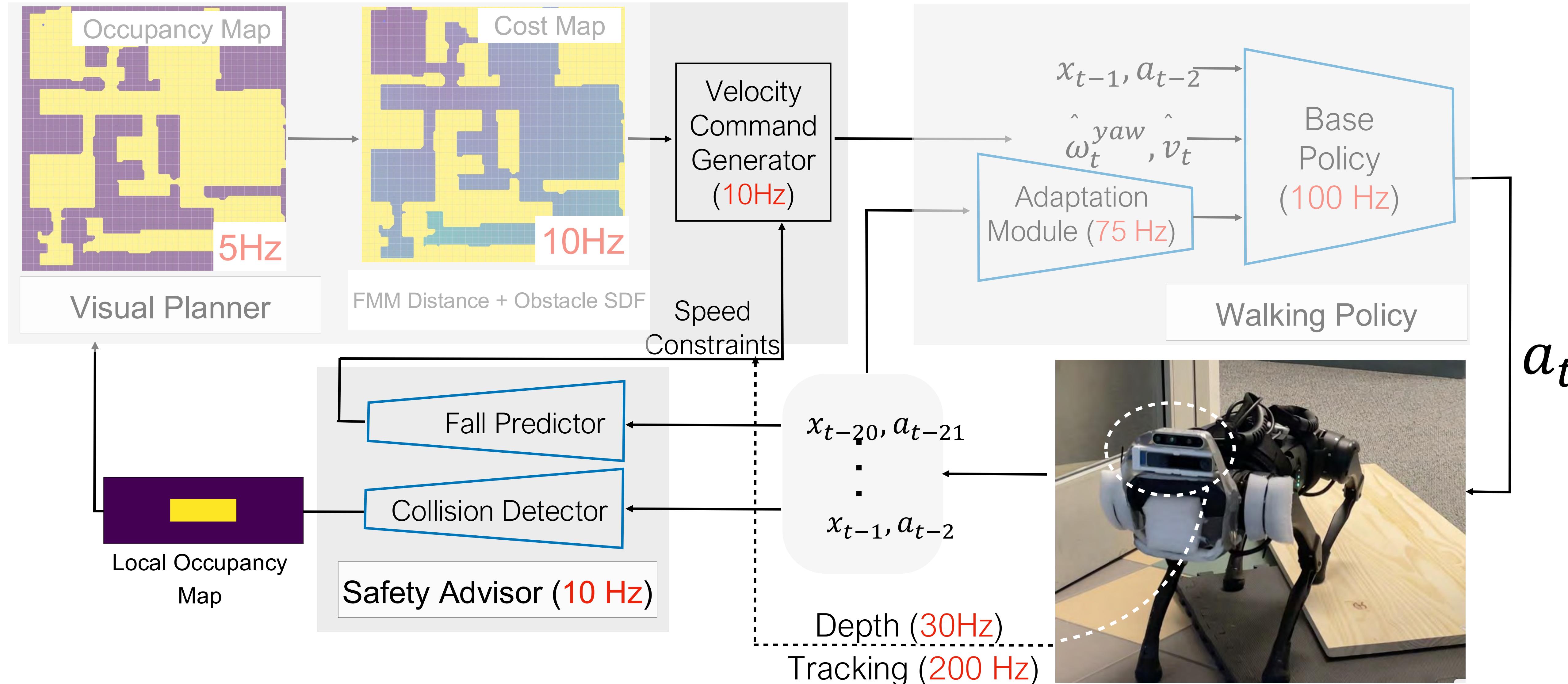
Vision coordinated with proprioception



Only Vision

Vision + Proprioception

Coupled Navigation and Locomotion

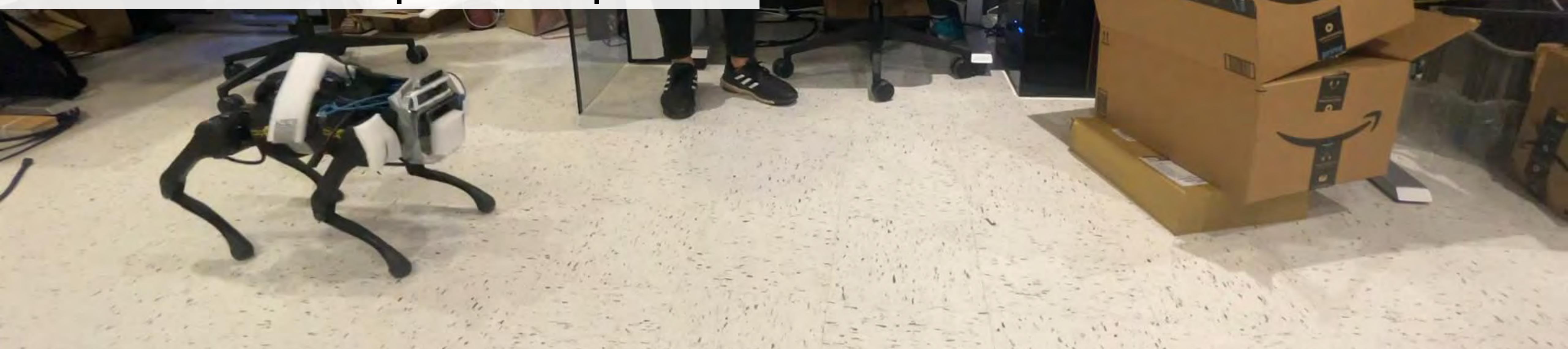


Comparison of

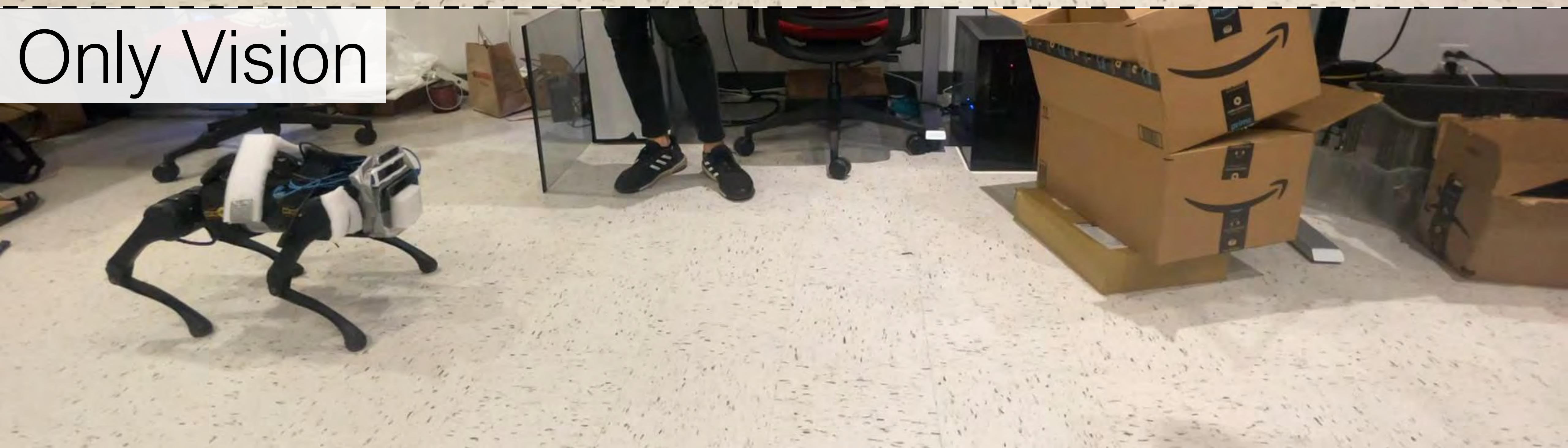
Vision + Proprioception

Only Vision

Vision + Proprioception



Only Vision



Vision + Proprioception

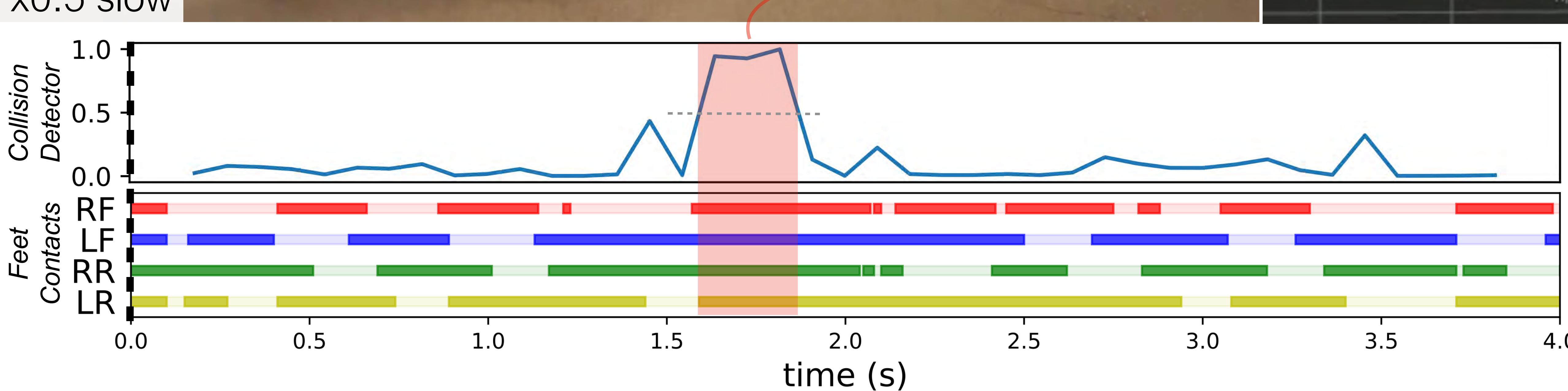


Only Vision



More Deployment Videos

Collision Detector

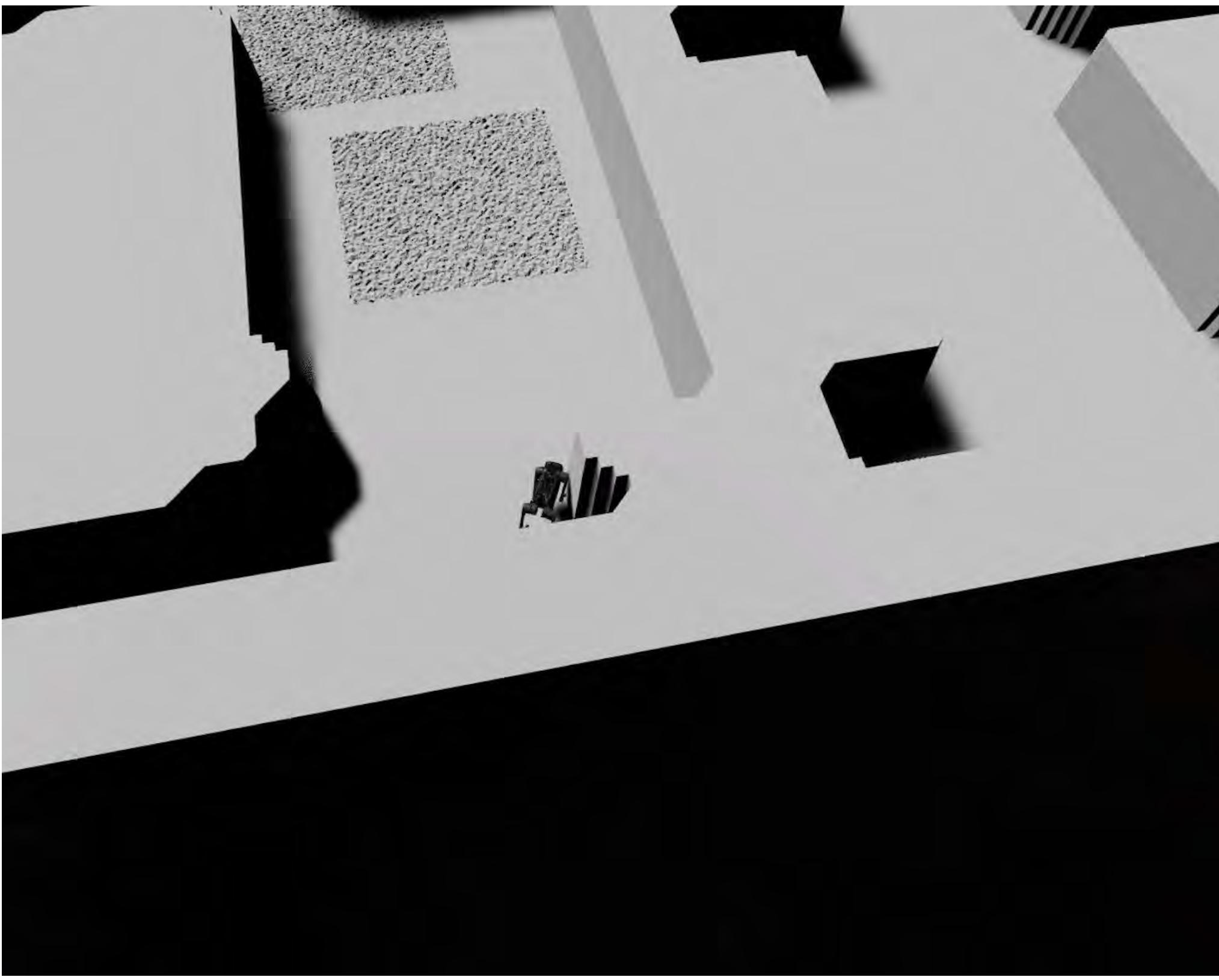


Navigation in Cluttered Indoor Setting



Navigation in the Wild





Simulation Comparisons

	Navigation System	Terrain Type	Success	SPL	Time (s)
(a)	w/o Proprio	Flat	95.20	0.79	80.28
(b)	w/o Proprio	2x Inv-Obstacle	68.45	0.57	119.80
(c)	Ours	2x Inv-Obstacle	74.15	0.61	111.93
(d)	w/o Proprio	4x Inv-Obstacle	45.85	0.38	152.39
(e)	Ours	4x Inv-Obstacle	59.20	0.49	134.70
(f)	w/o Proprio	8x Inv-Obstacle	24.35	0.20	184.07
(g)	Ours	8x Inv-Obstacle	39.25	0.32	164.95

invisible obstacles

	Navigation System	Terrain Type	Success	SPL	Time (s)
(a)	w/o Proprio	Flat	95.20	0.79	80.28
(b)	w/o Proprio	Randomized	80.25	0.66	105.68
(c)	Ours	Randomized	87.40	0.73	117.65

rough, slip, payload randomization

Real world comparisons

Detour with Planks

	Success Rate	# Obstacles Hit
Vision + Proprio	100%	0.2
Vision	80%	0.2

Cluttered Narrow Path

	Success Rate	# Obstacles Hit
Vision + Proprio	80%	0.6
Vision	60%	0.6

Unexpected Human Obstacle

	Success Rate	Time to Recover After Blocked (s)
Vision + Proprio	100%	0.9
Vision	0%	∞

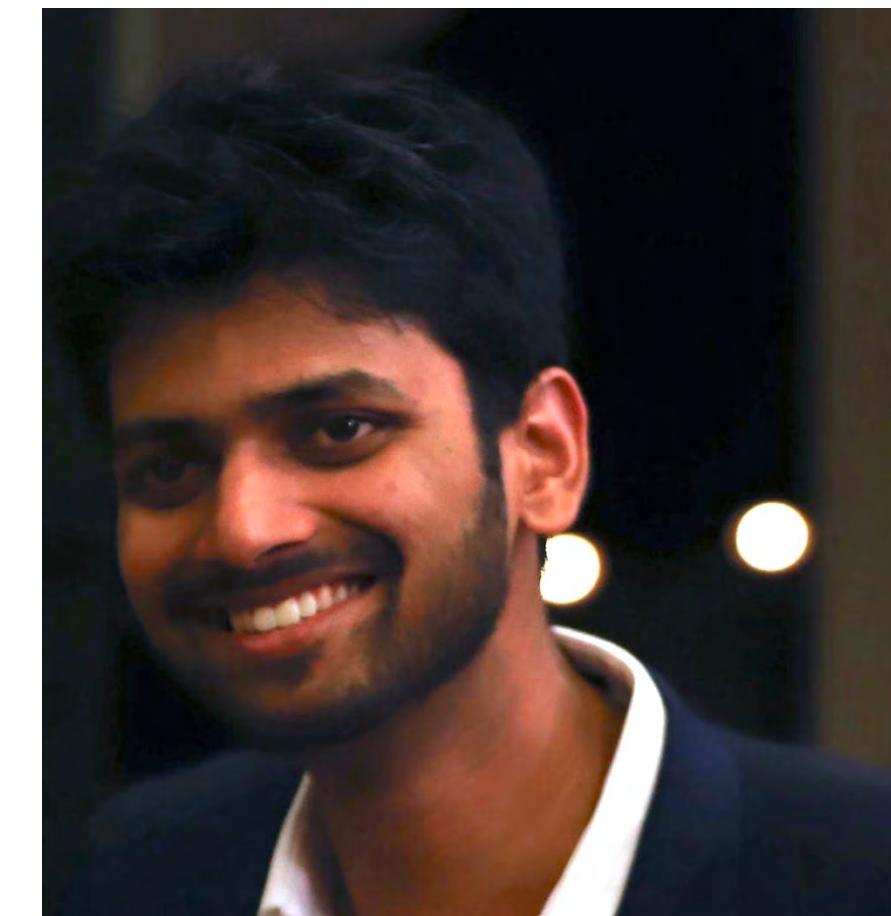
Rough Slippery Terrain

	Success Rate	# Slips
Vision + Proprio	100%	3.2
Vision	80%	4.0

Legged Locomotion in Challenging Terrains using Egocentric Vision



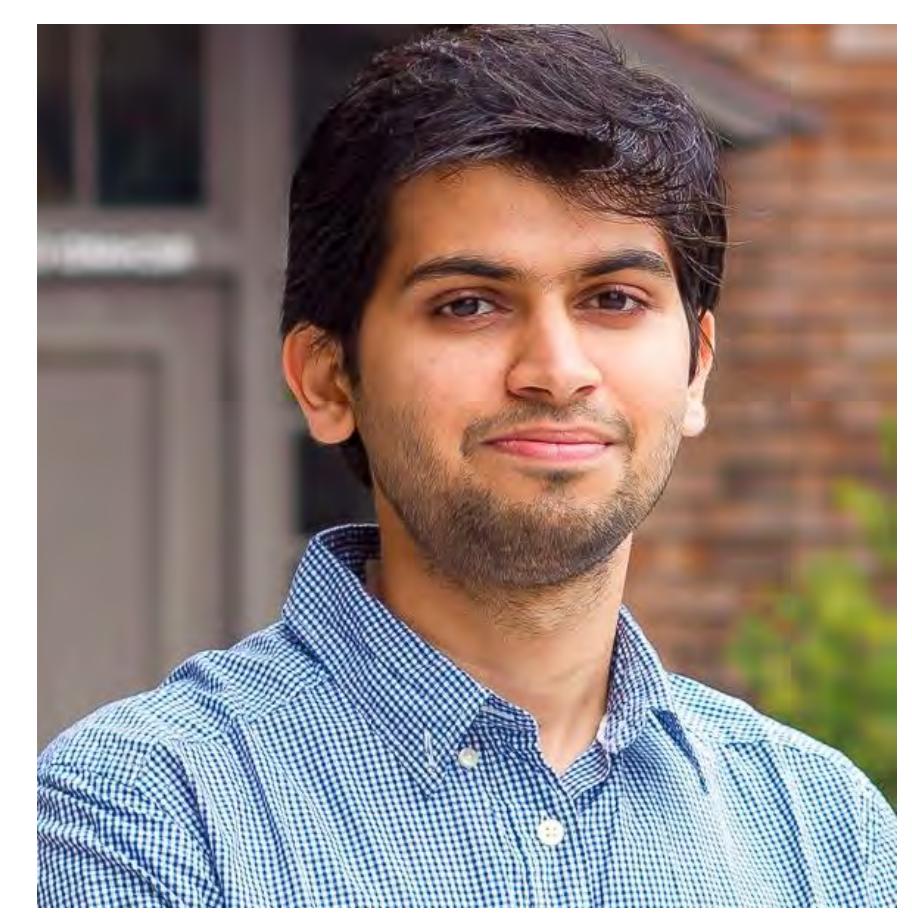
Ananye Agarwal*
CMU



Ashish Kumar*
UC Berkeley



Jitendra Malik†
UC Berkeley



Deepak Pathak†
CMU

Robots can walk on challenging terrain

These robots are blind!



Kumar, Ashish, et al.
"Rma: Rapid motor
adaptation for legged
robots." RSS (2021).

Siekmann, Jonah, et al. "Blind
bipedal stair traversal via sim-to-
real reinforcement
learning." RSS (2021)

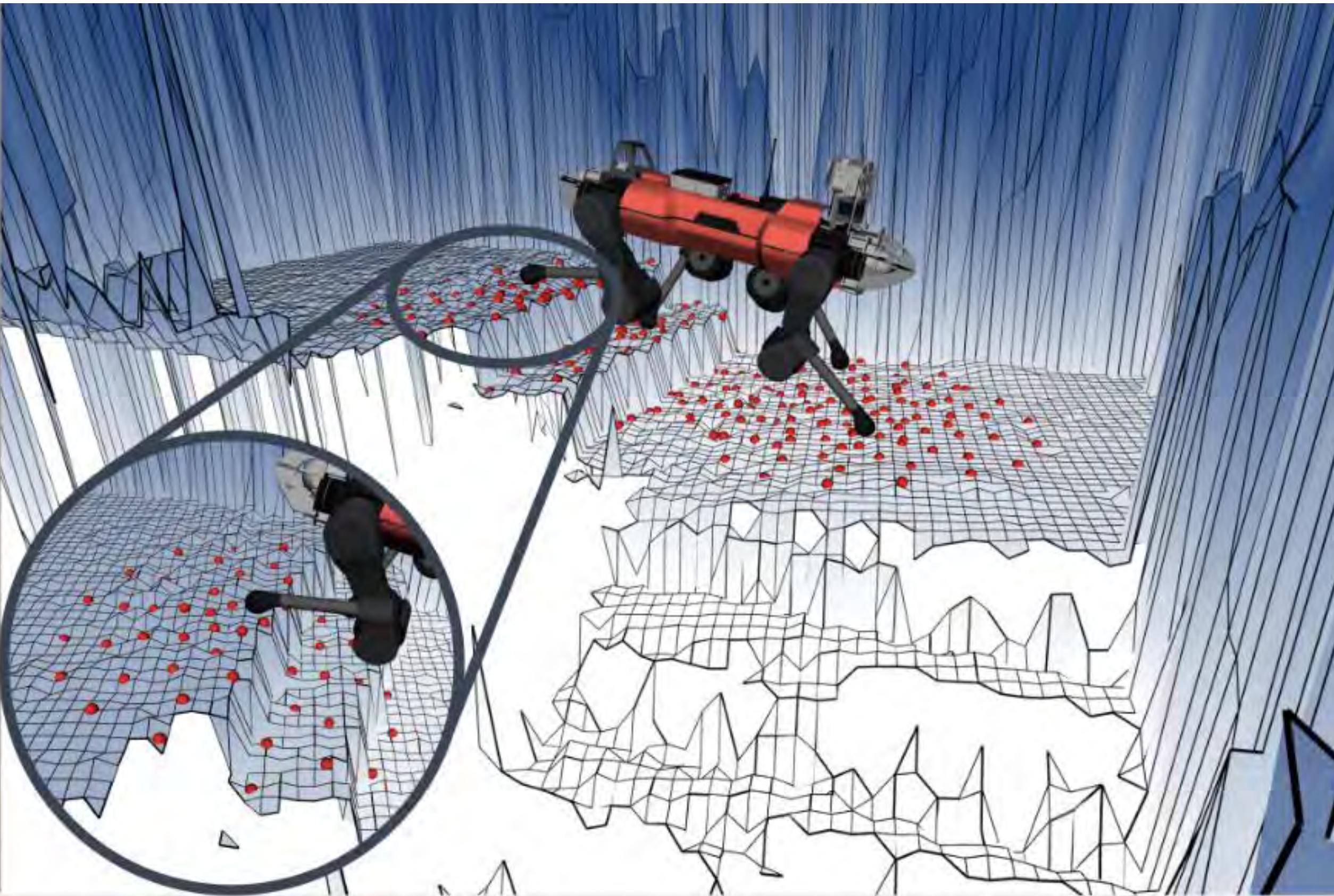
Lee, Joonho, et al. "Learning
quadrupedal locomotion over
challenging terrain." Science
robotics 5.47 (2020)

Why do we need vision?

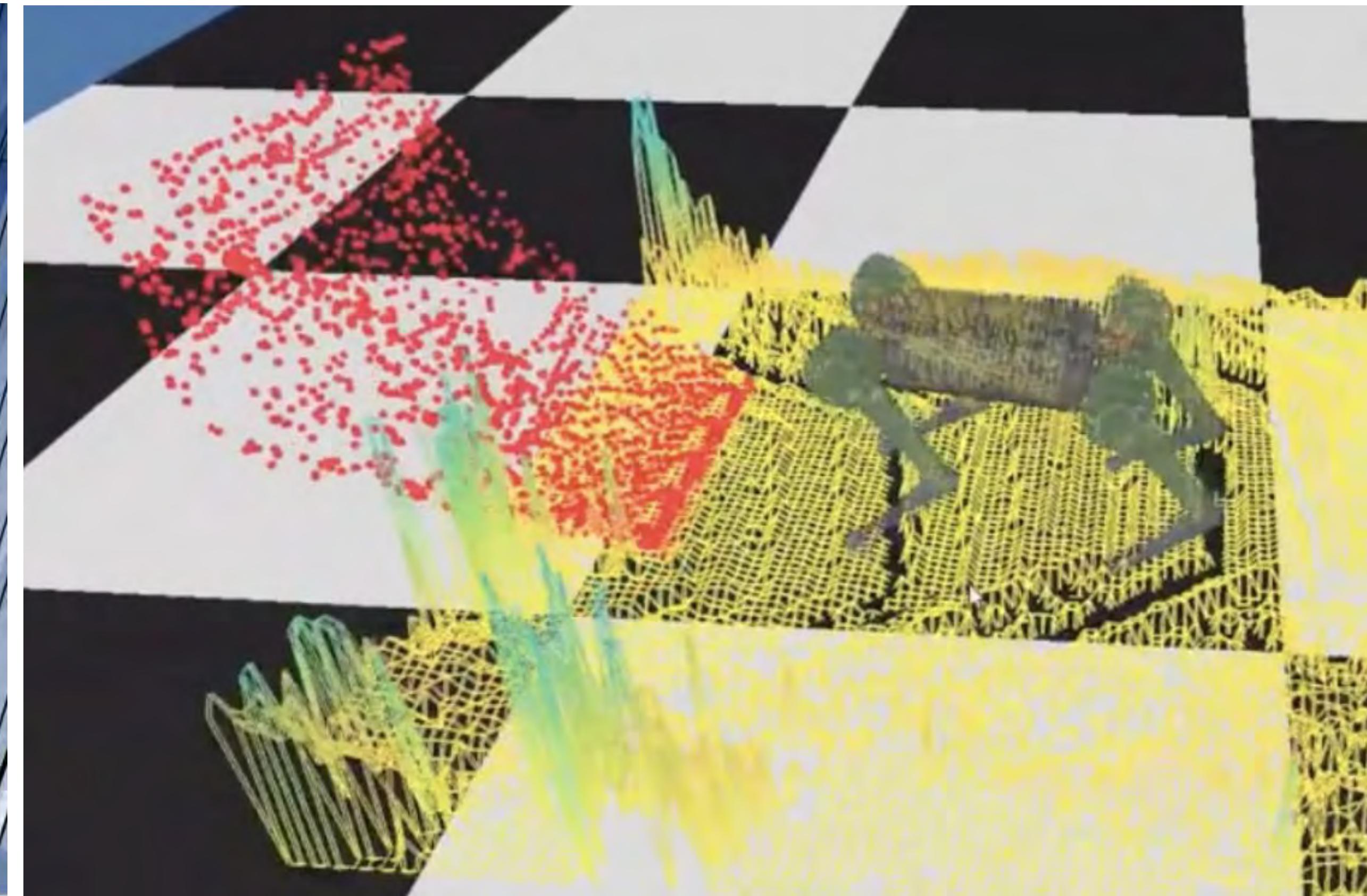
Perception enables precise locomotion



Typical approach: build terrain maps from vision



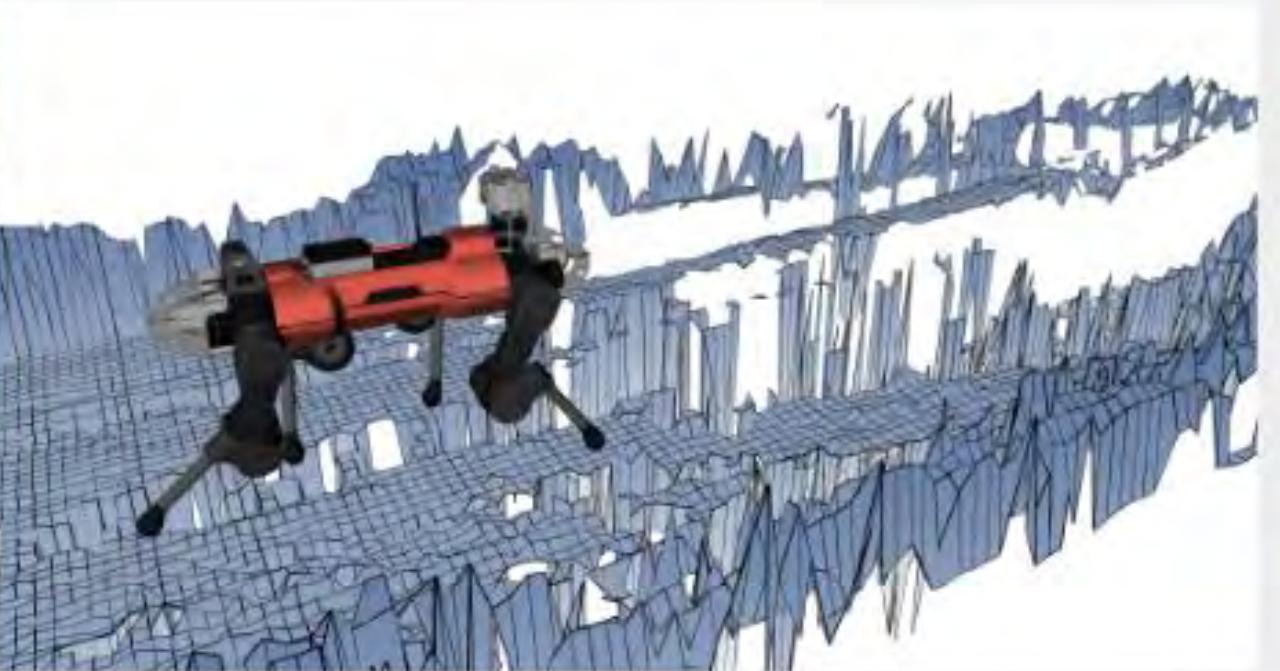
Miki, T., Lee, J., Hwangbo, J., Wellhausen, L., Koltun, V., & Hutter, M. (2022). Learning robust perceptive locomotion for quadrupedal robots in the wild. *Science Robotics*



Kim, Donghyun, et al. "Vision aided dynamic exploration of unstructured terrain with a small-scale quadruped robot." ICRA 2020.

Real world maps accumulate noise

B Reflective ground



E Non rigid obstacles



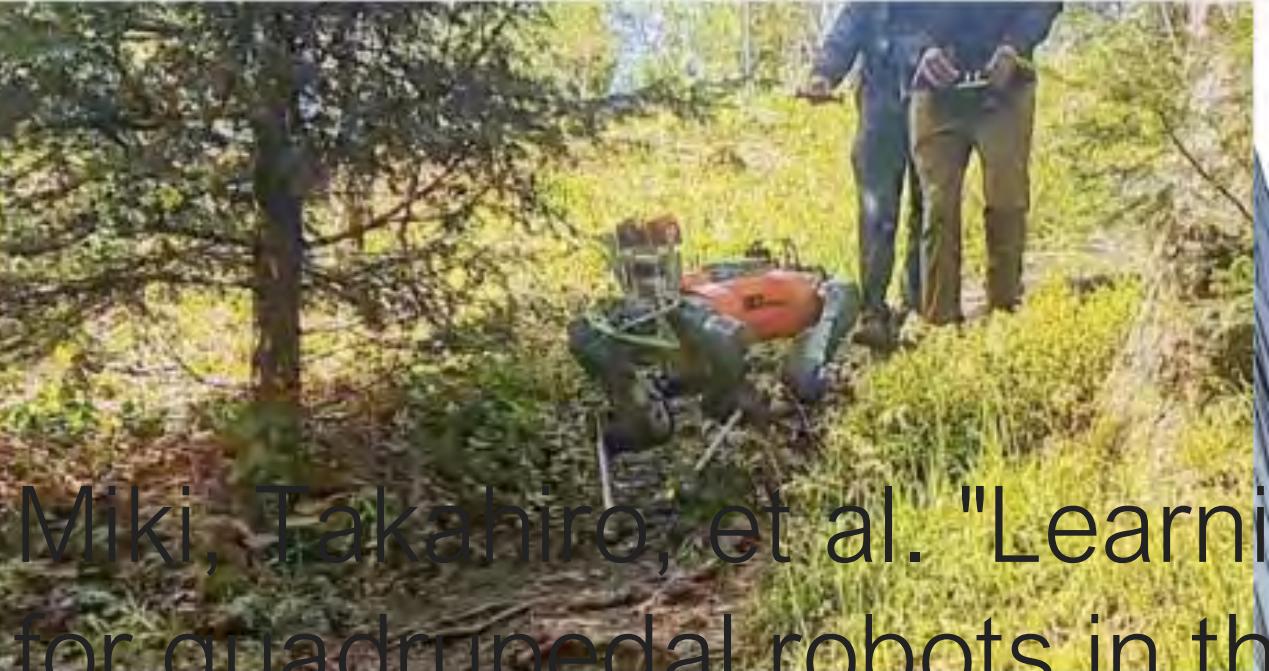
C Deep snow



F Pose estimation drift



D Overhanging objects



G Occlusion



Miki, Takahiro; et al. "Learning robust perceptive locomotion for quadrupedal robots in the wild." *Science Robotics* (2022)

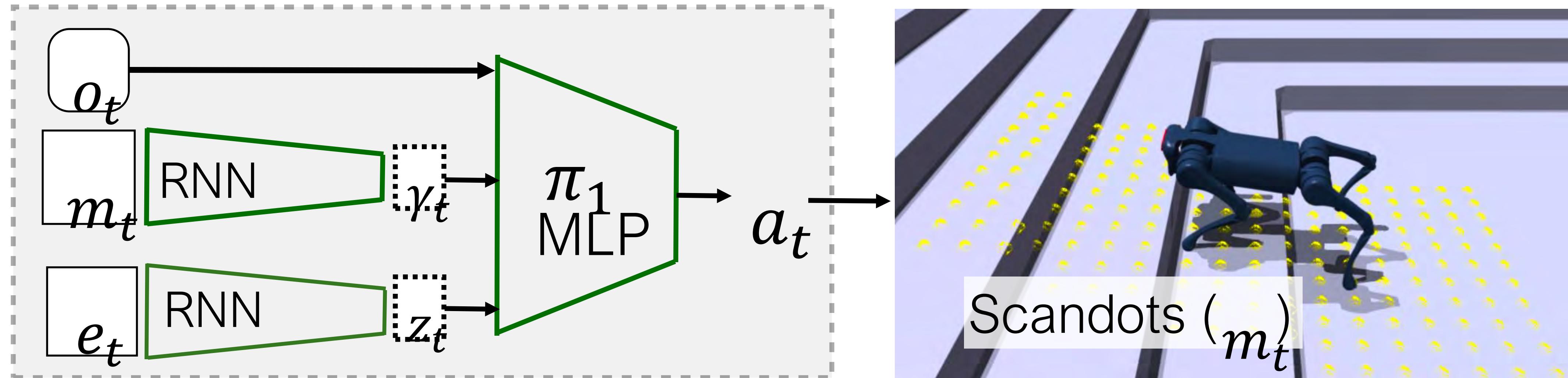
Do we really need terrain maps?

We directly go from vision to control

Stepping Stones (~15 cm apart)

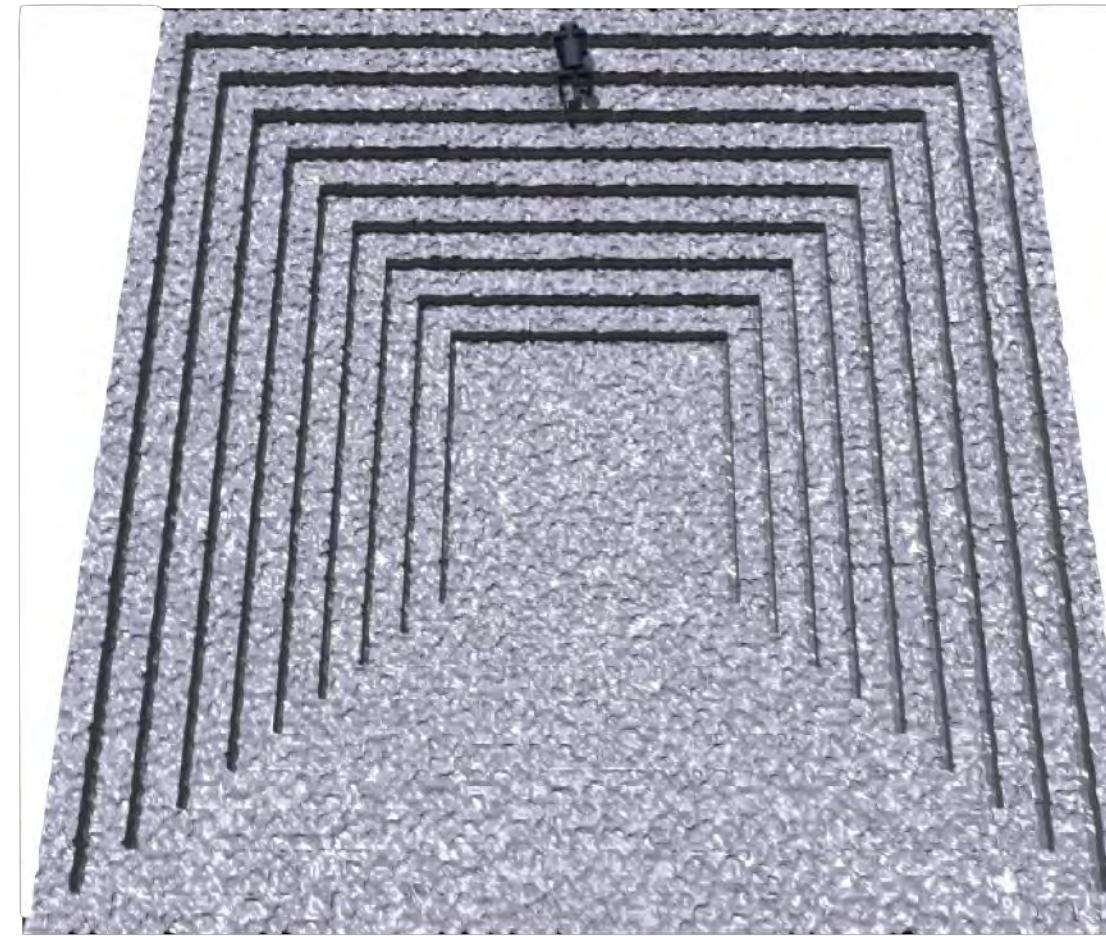


Phase 1: Learning to Walk with Privileged Terrain Information

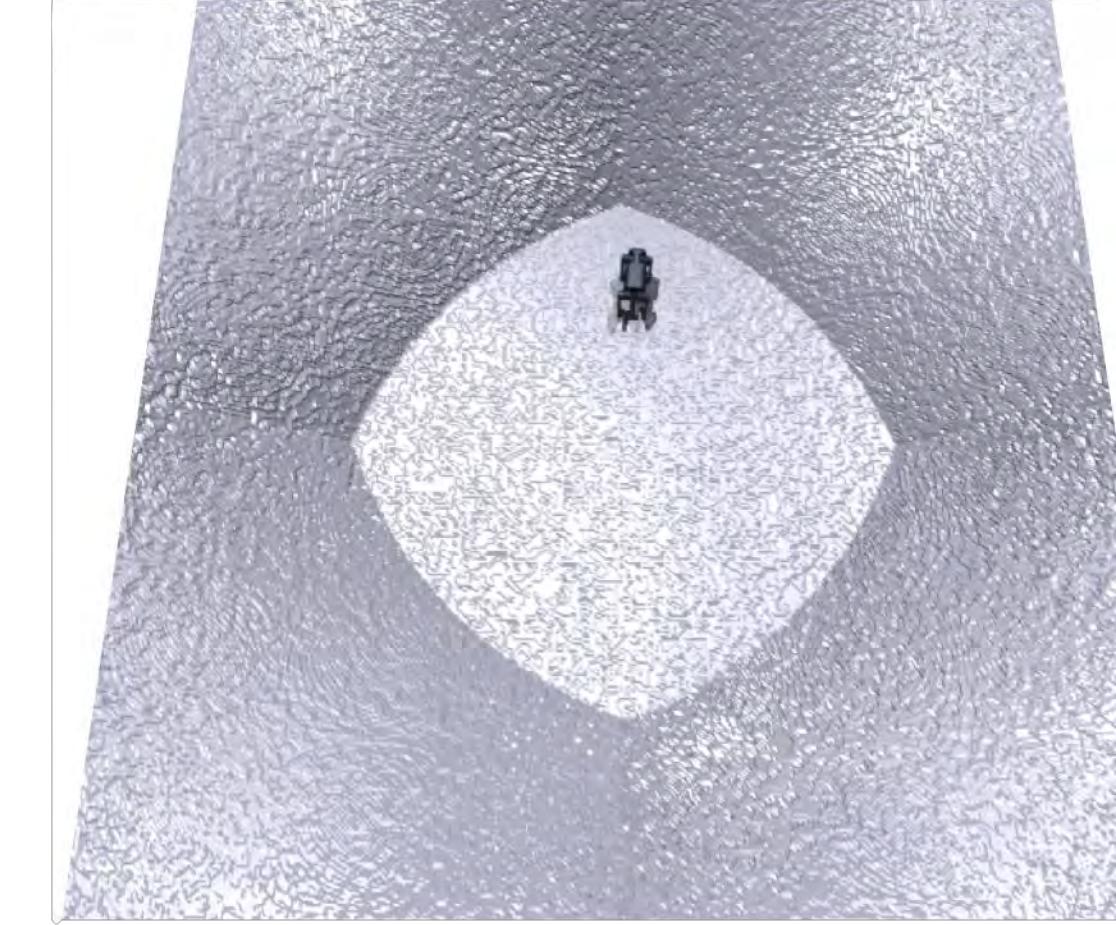


- Simulation: IsaacGym
- Trained with PPO with smoothness+task rewards
- ~10 Billion Samples simulated in 3 days on a single GPU

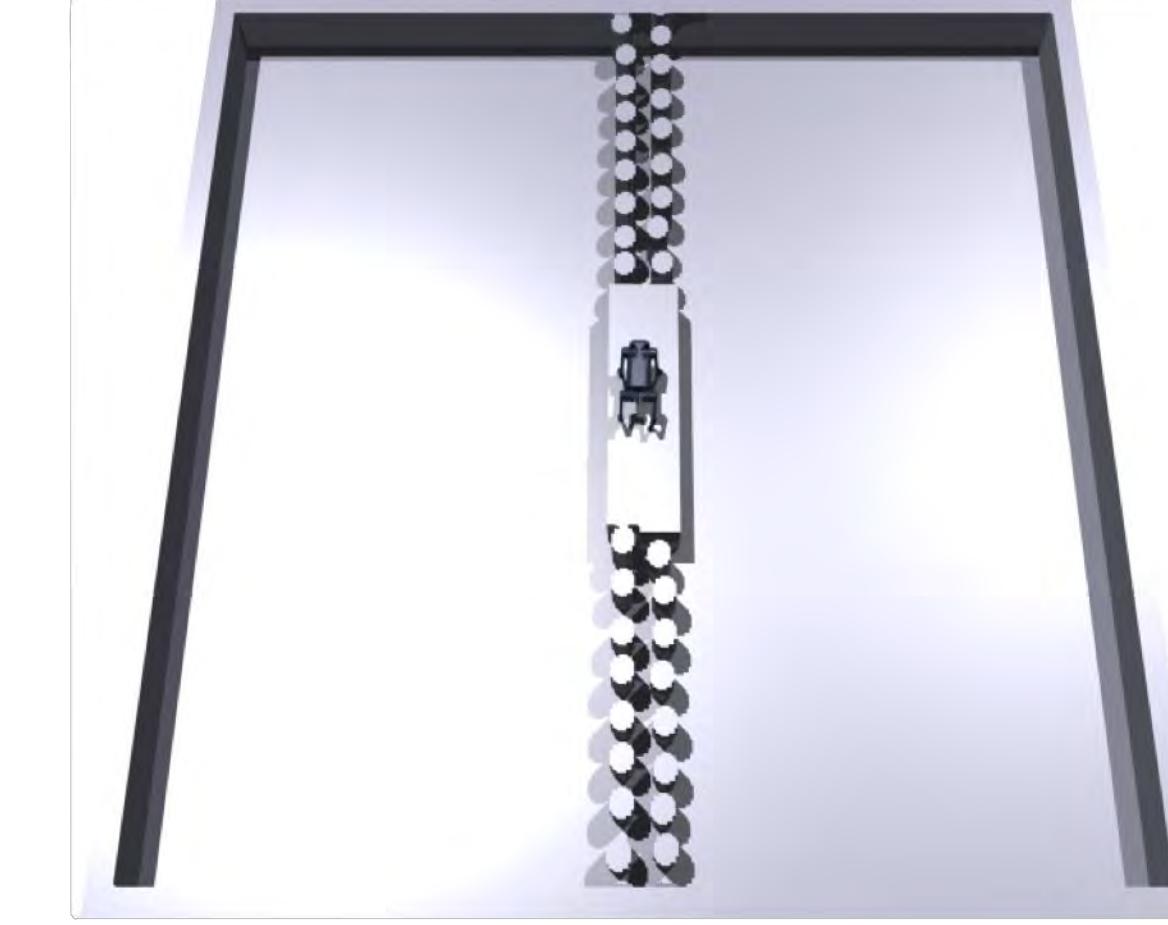
Training Terrains



Stairs



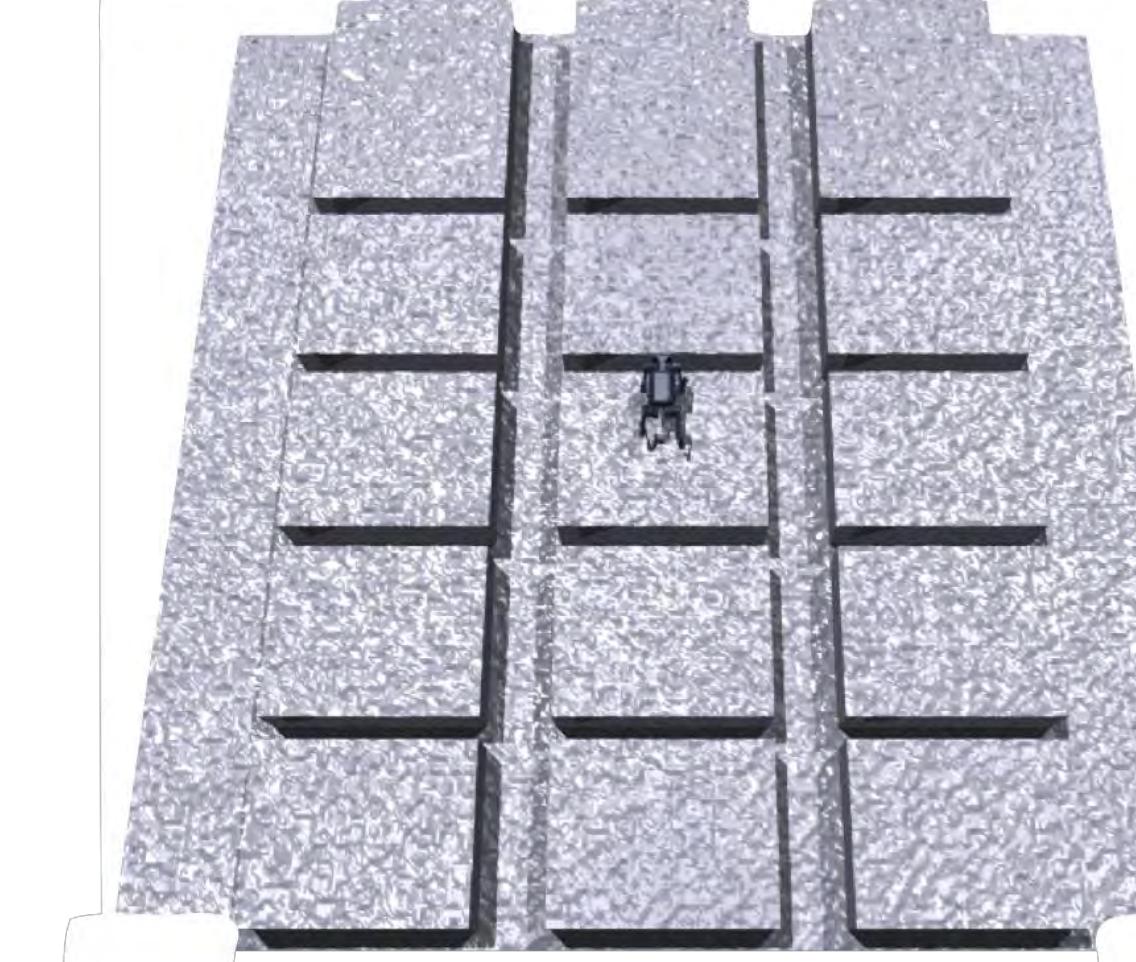
Slopes



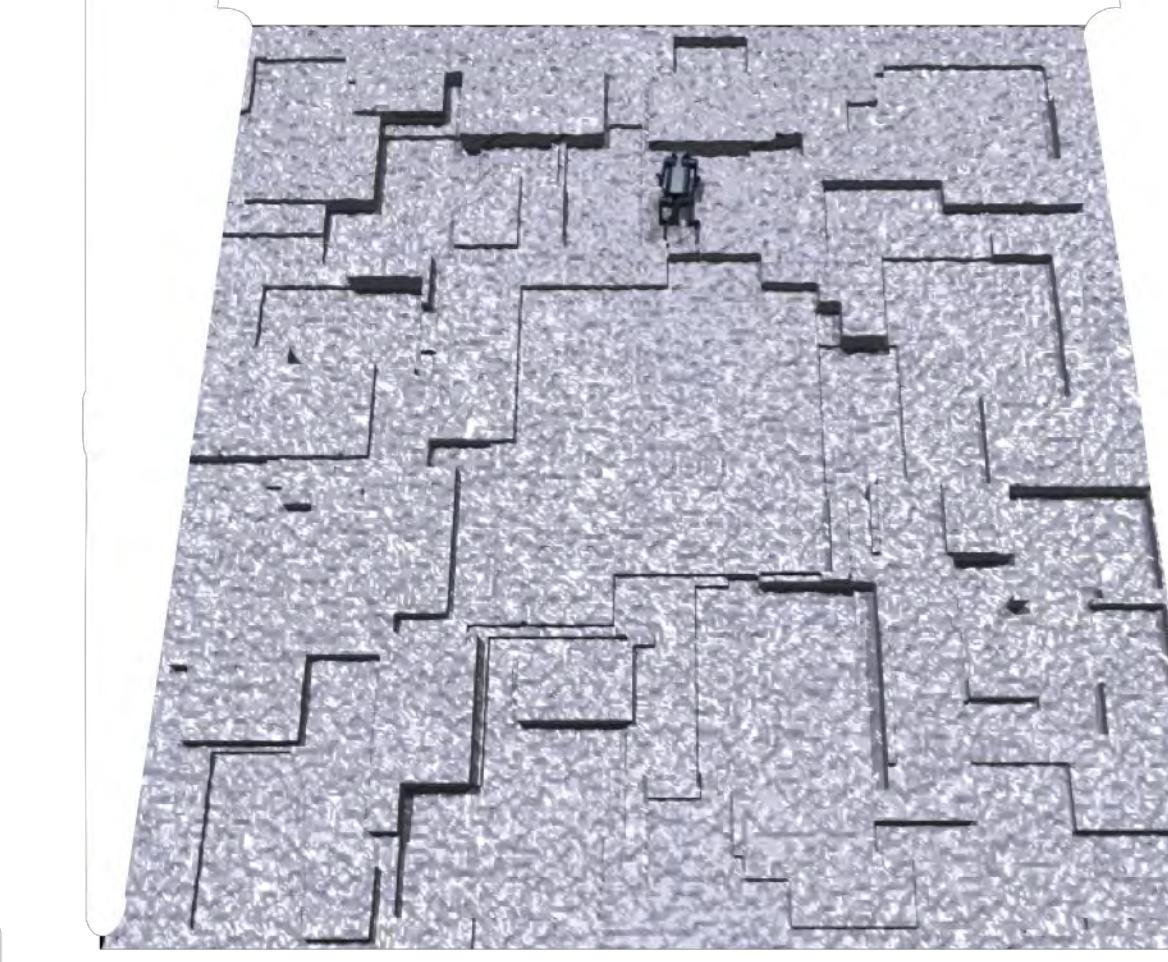
Stepping Stones



Rough Flat

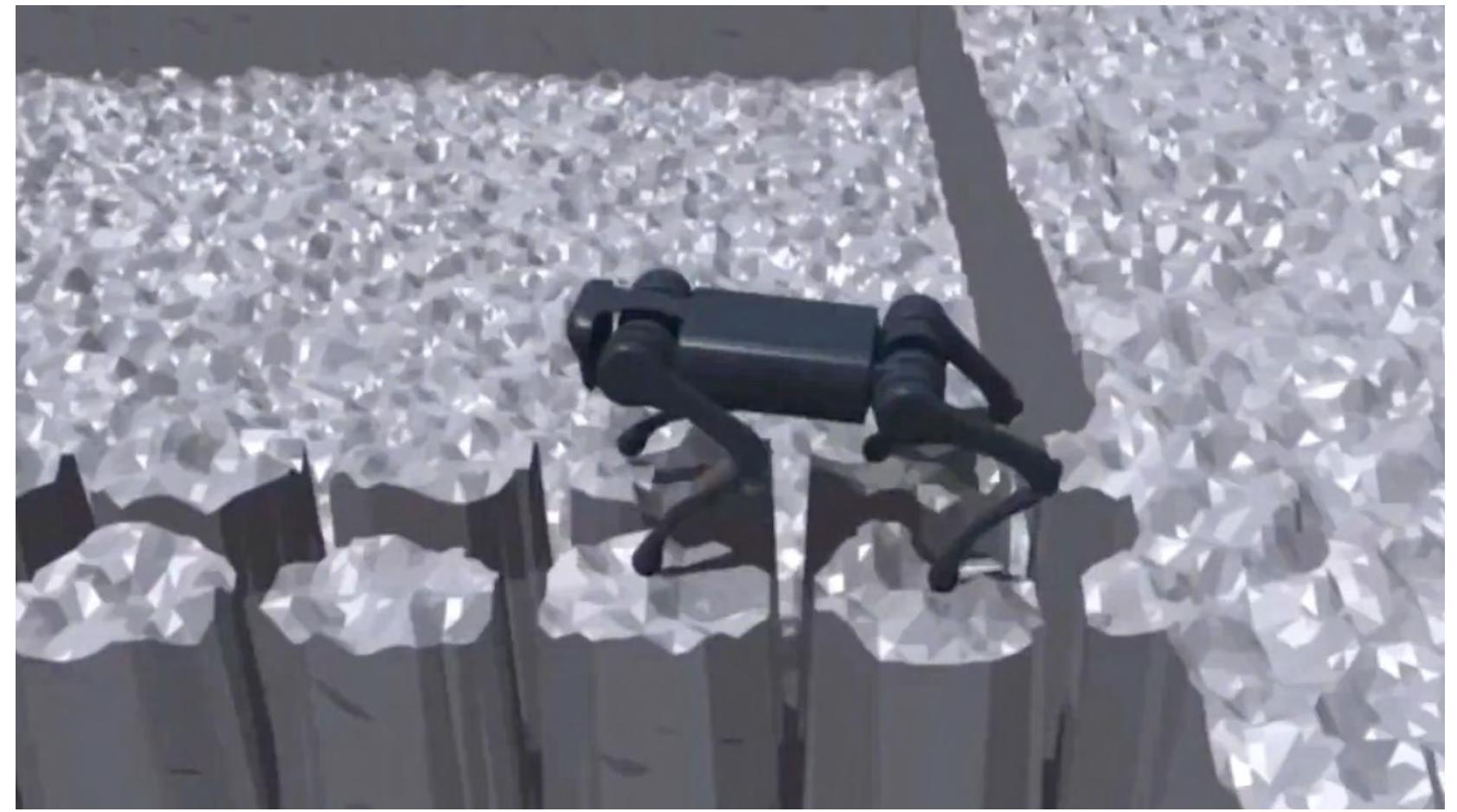


Gaps

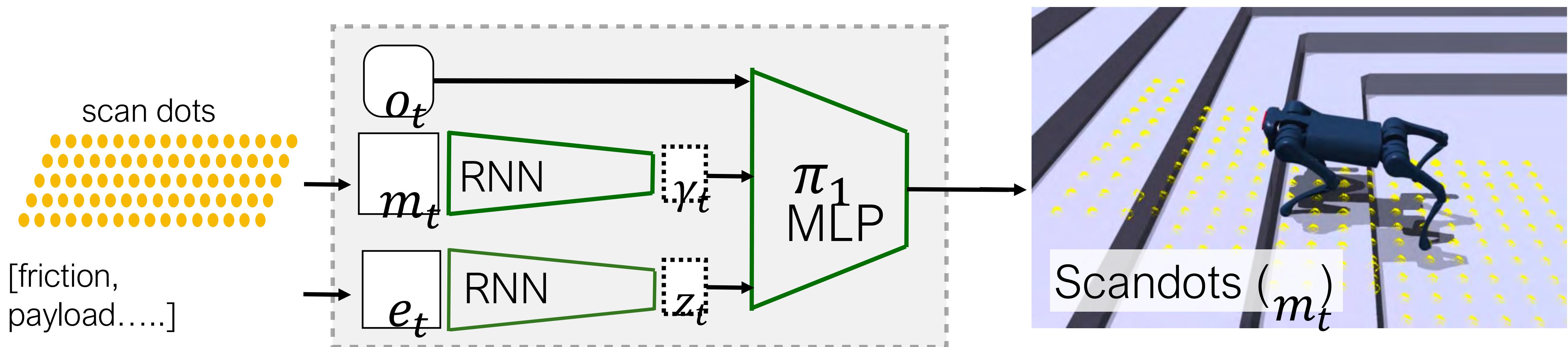


Discrete Obstacles

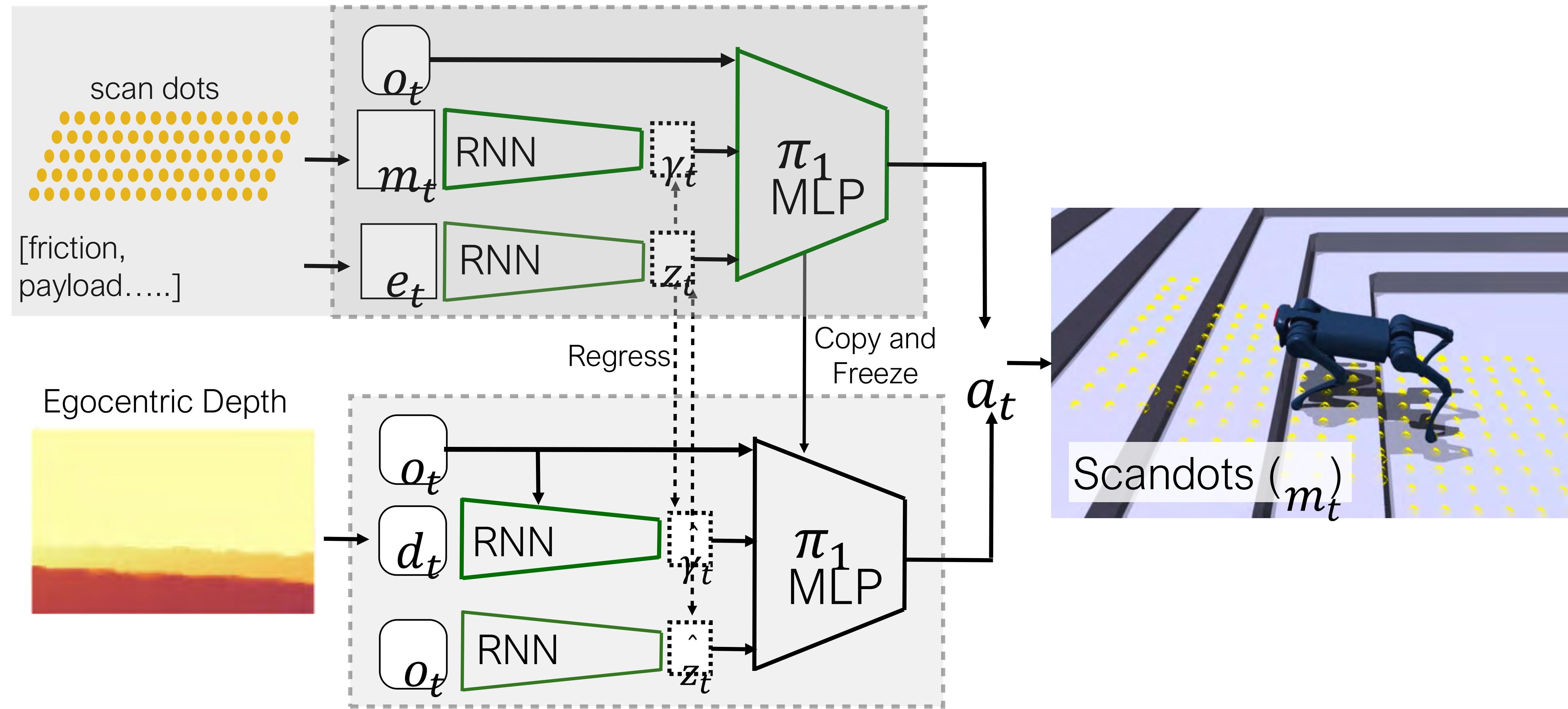
Phase 1 Policy



How do we deploy it?

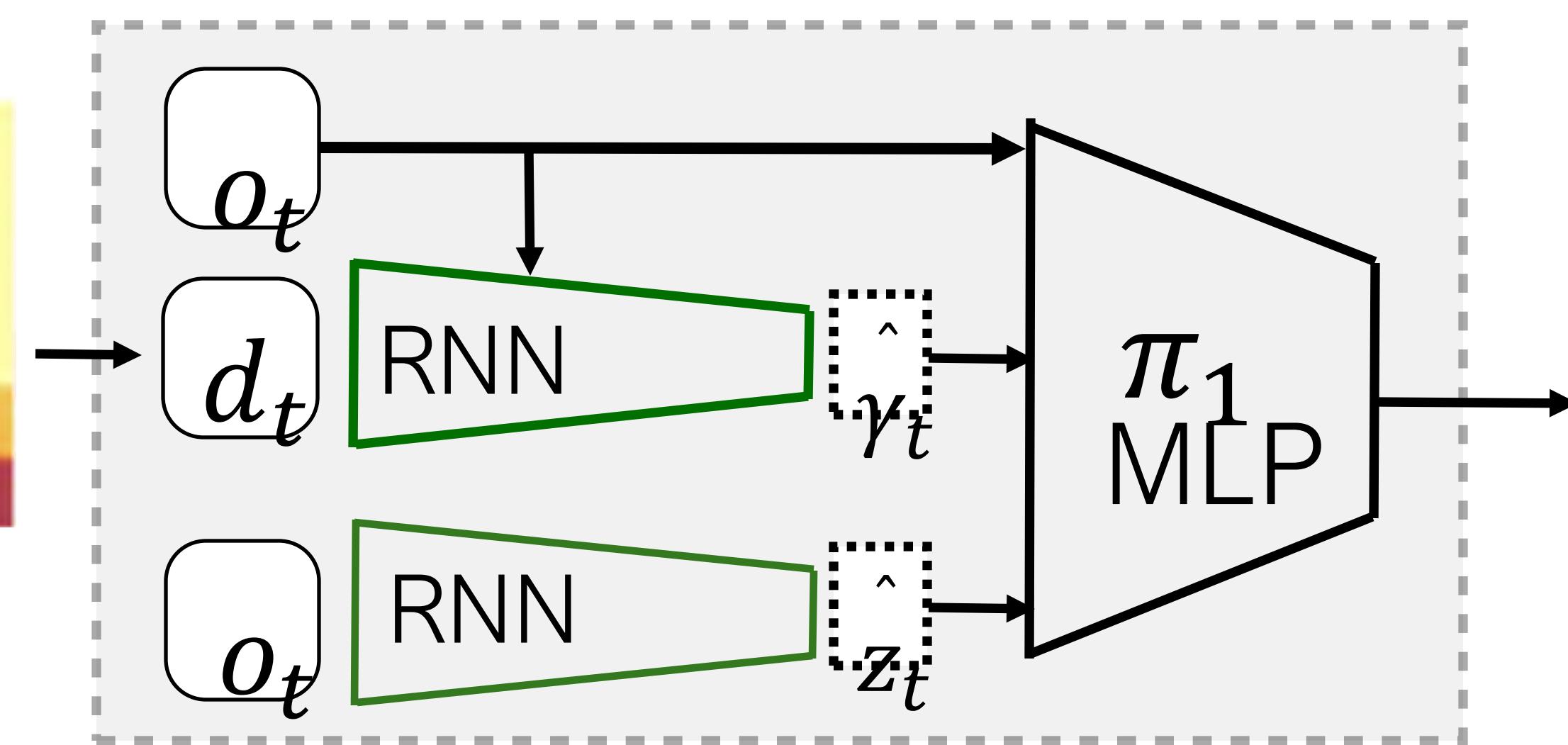


Phase 2: Learning to Walk with Egocentric Depth



Deployment Policy

Egocentric Depth

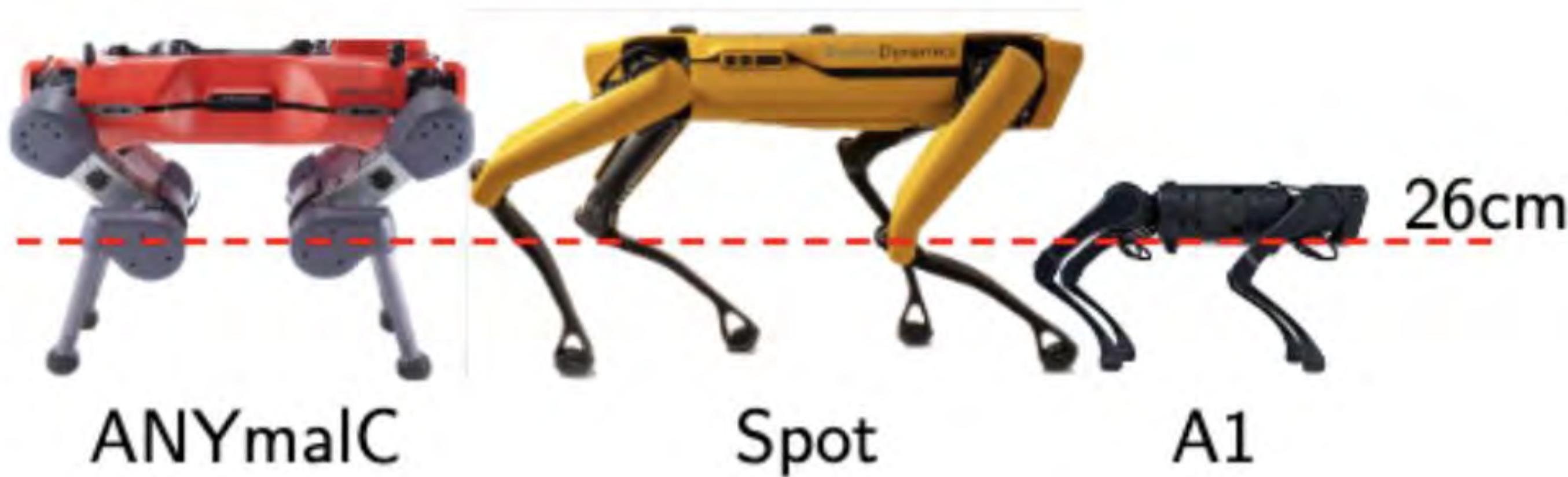


Stepping Stones (~20 cm apart)

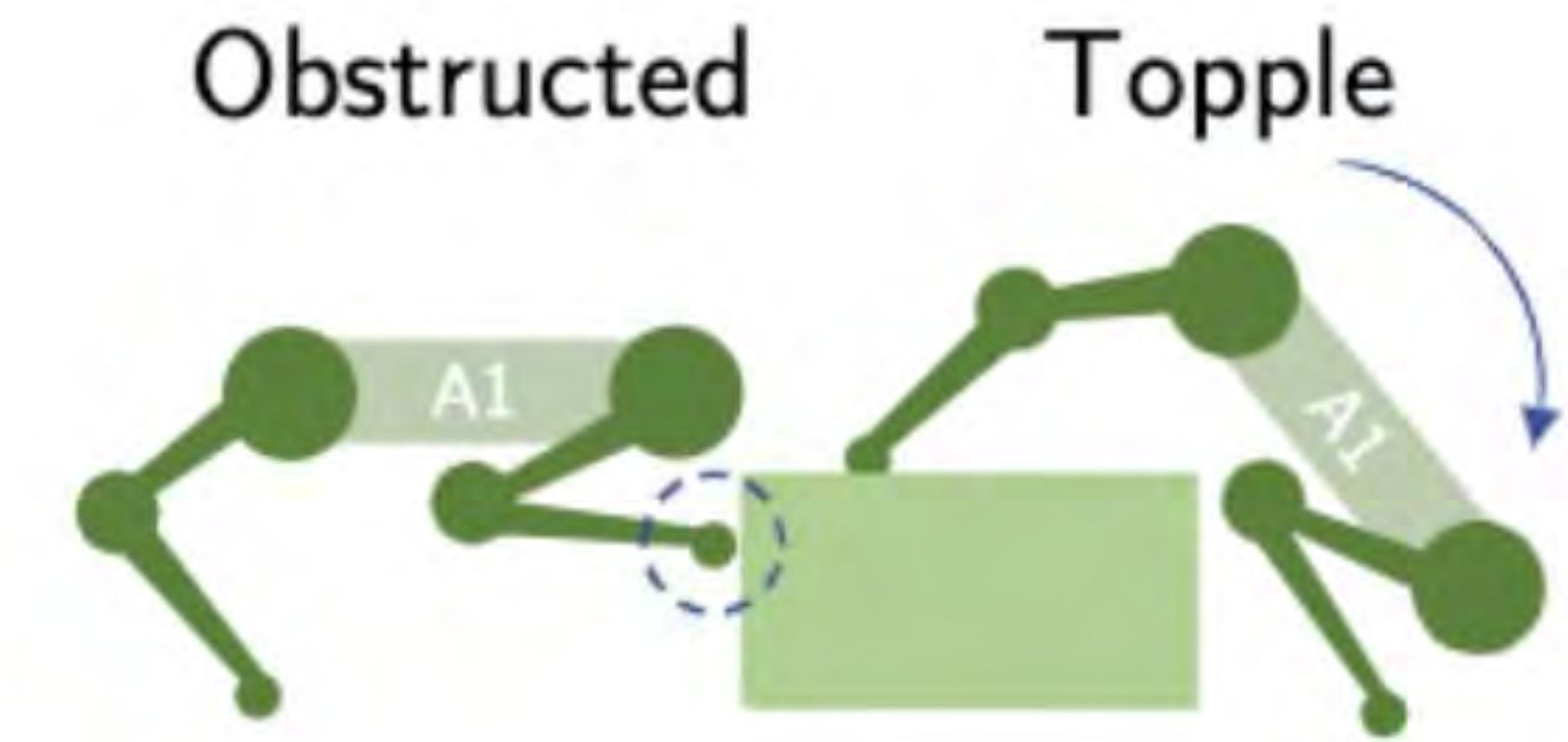


Emergent footstep planning

Gait heuristics fail on a small robot



Size Comparison

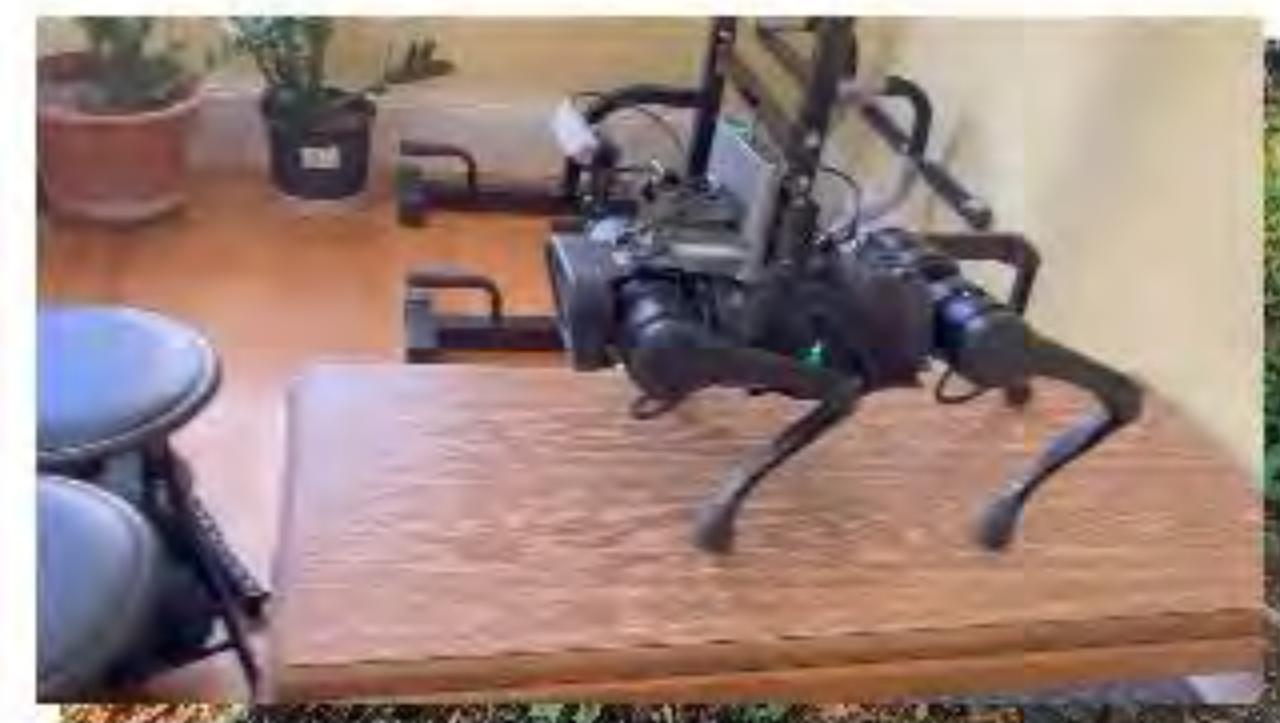


Challenges Due to Size

Emergent Hip Abduction



Map Free, Gait Free



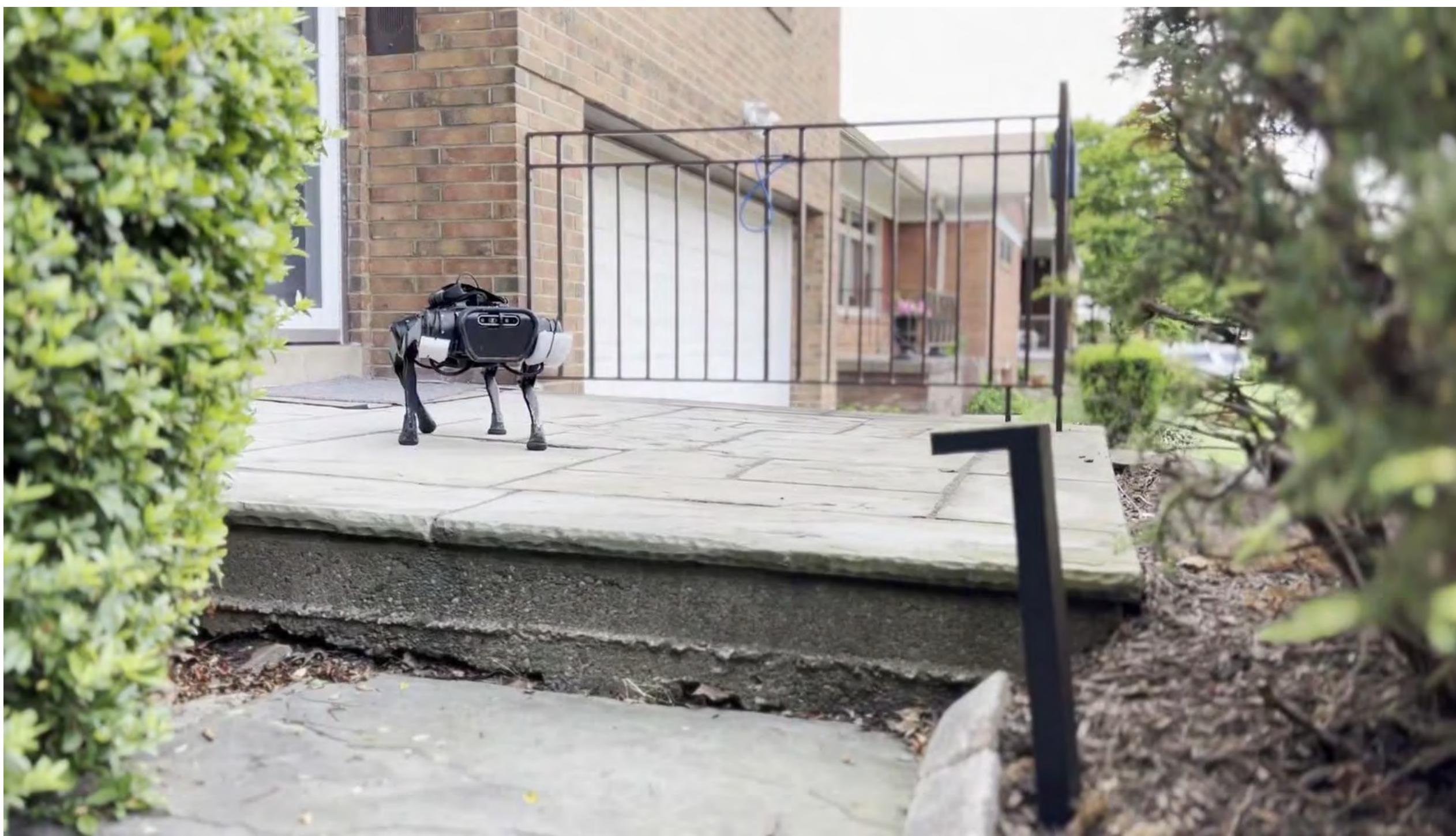
Performance is better without maps

Average Distance

	Blind	Noisy	Ours
Slopes	34.72	36.14	43.98
Stepping Stones	1.02	1.09	18.83
Stairs	16.64	6.74	31.24
Discrete Obstacles	32.41	29.08	40.13

Limitations

Front Camera Failure (unable to see a dip)



Implicit planning failure





Thank You