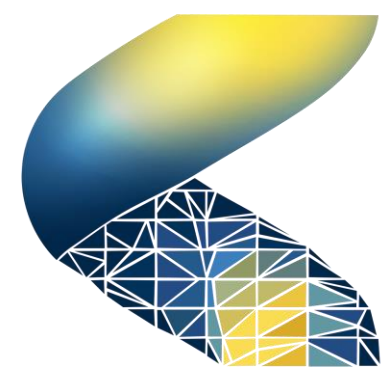


# A Viewpoint on Perception, Planning, and Control

**Claire Tomlin**

**Somil Bansal, Varun Tolani, Aleksandra Faust, Jitendra Malik**

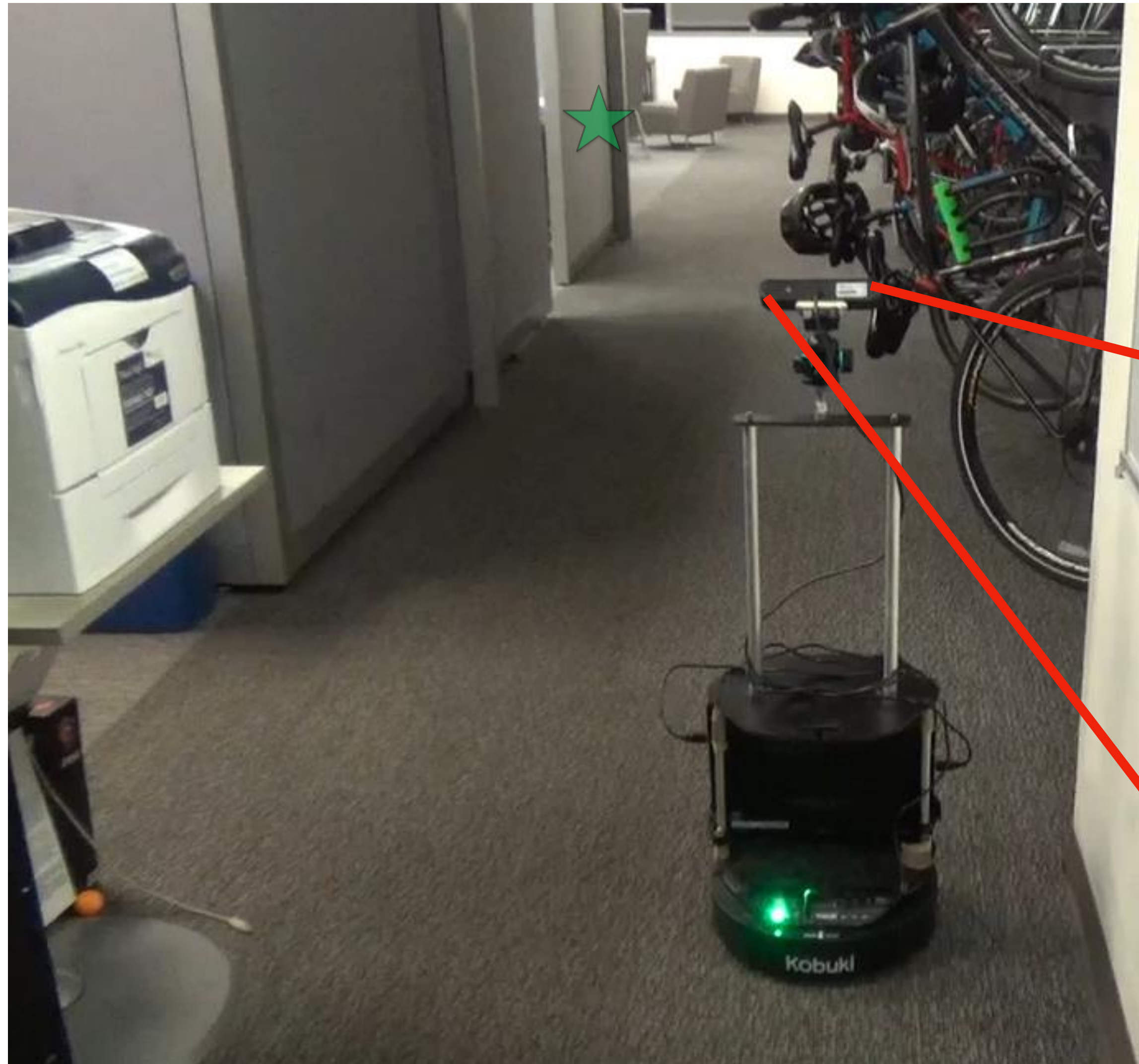


Hybrid Systems  
Laboratory



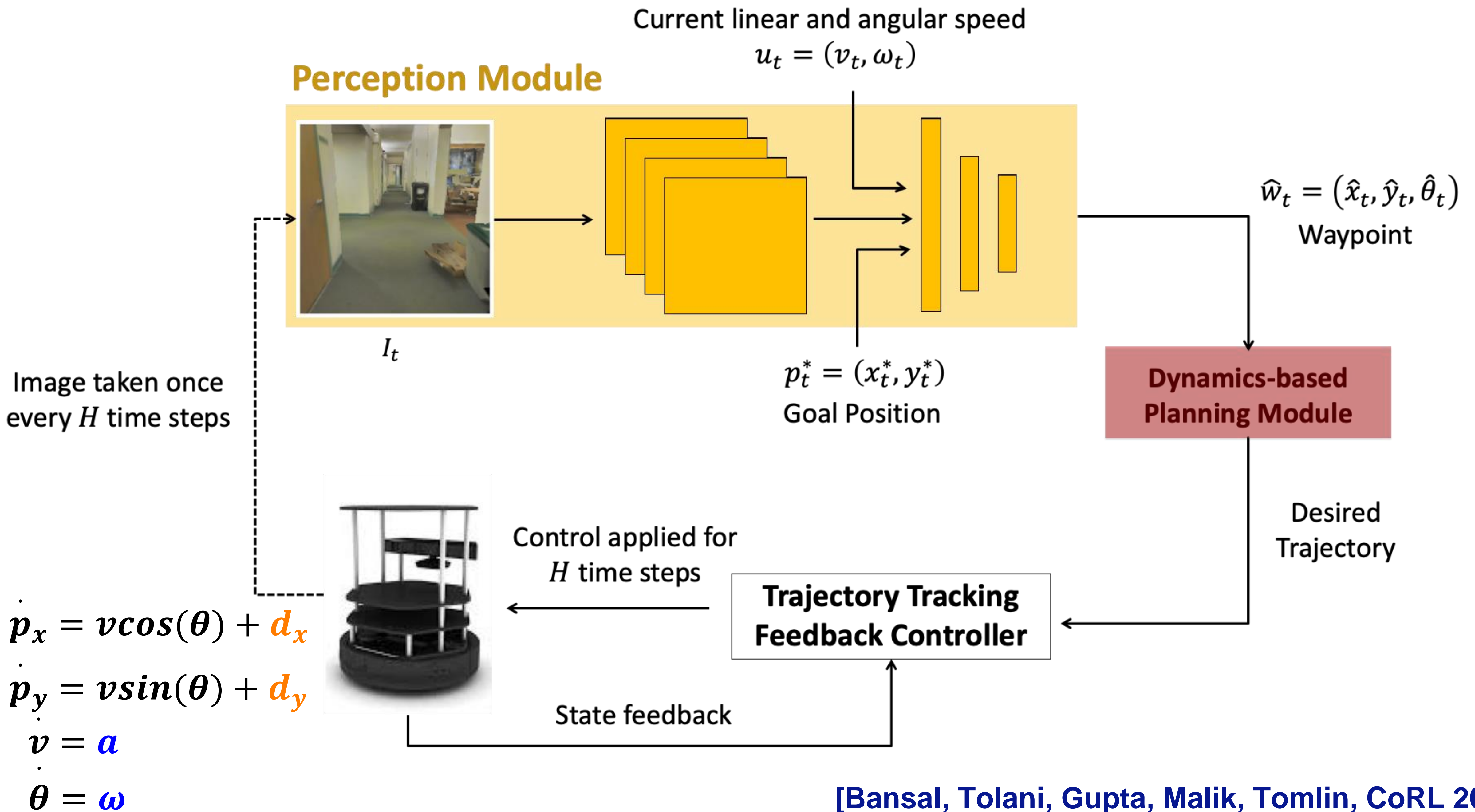


How to efficiently navigate an autonomous system with a monocular RGB camera to a goal in an *a priori* unknown environment?



[Bansal, Tolani, Gupta, Malik, Tomlin, CoRL 2019]





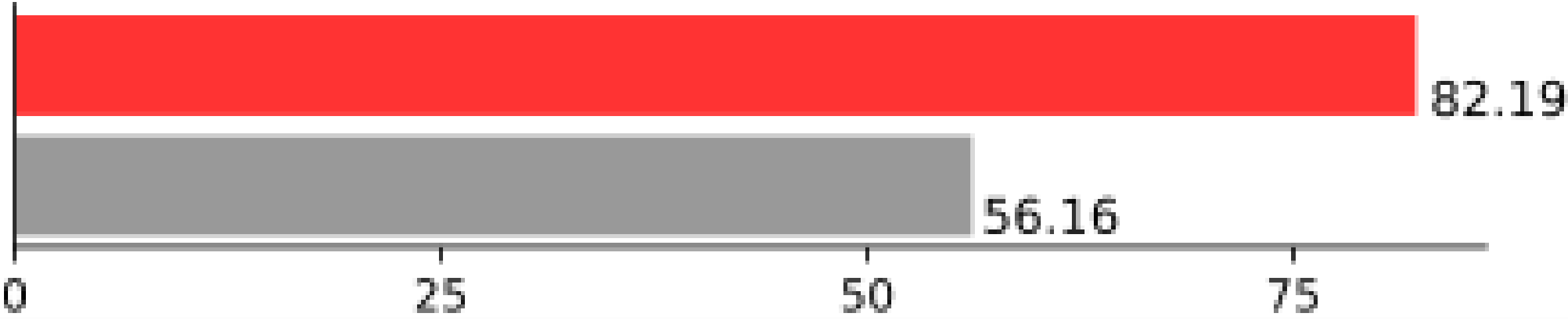


model-based

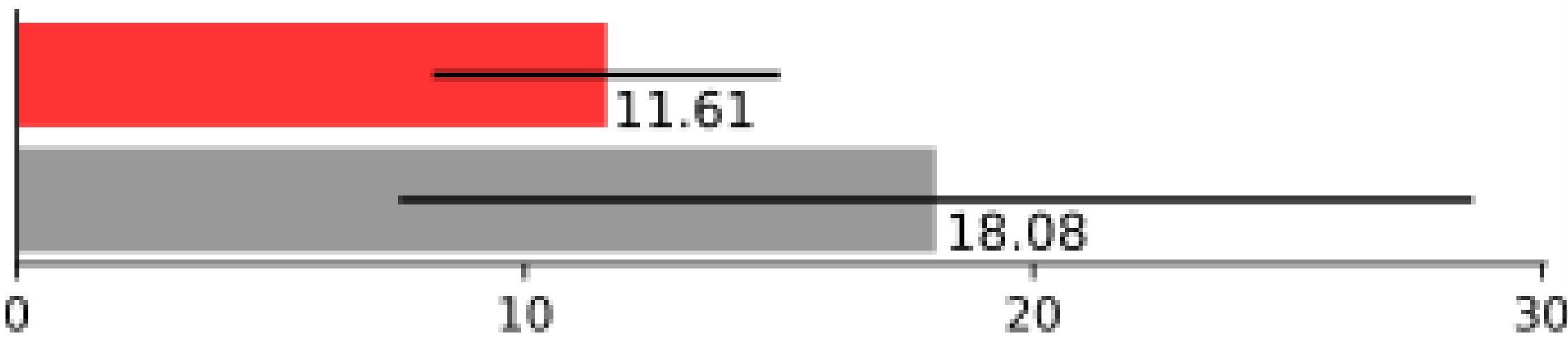


E2E

Success rate in reaching the goal (%):



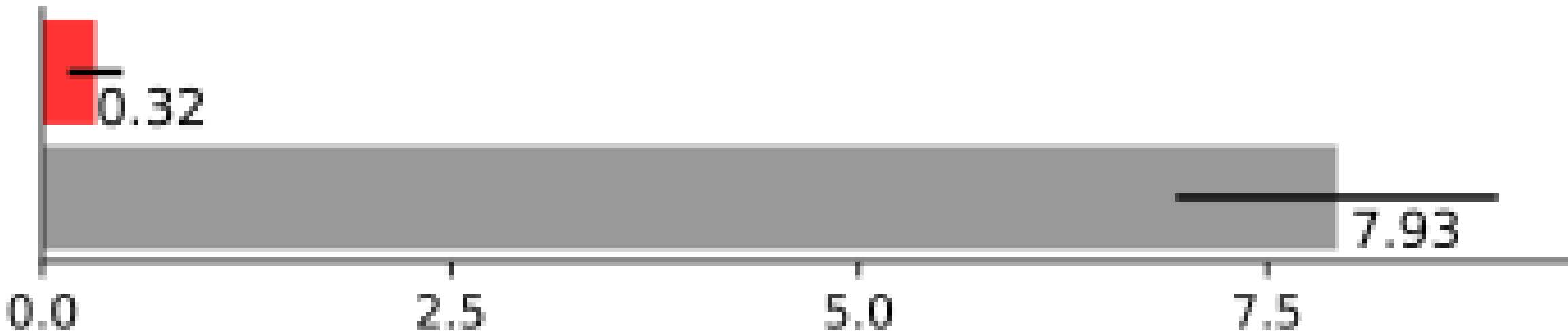
Time taken to reach the goal (s):



Average acceleration along the trajectory ( $m/s^2$ ):

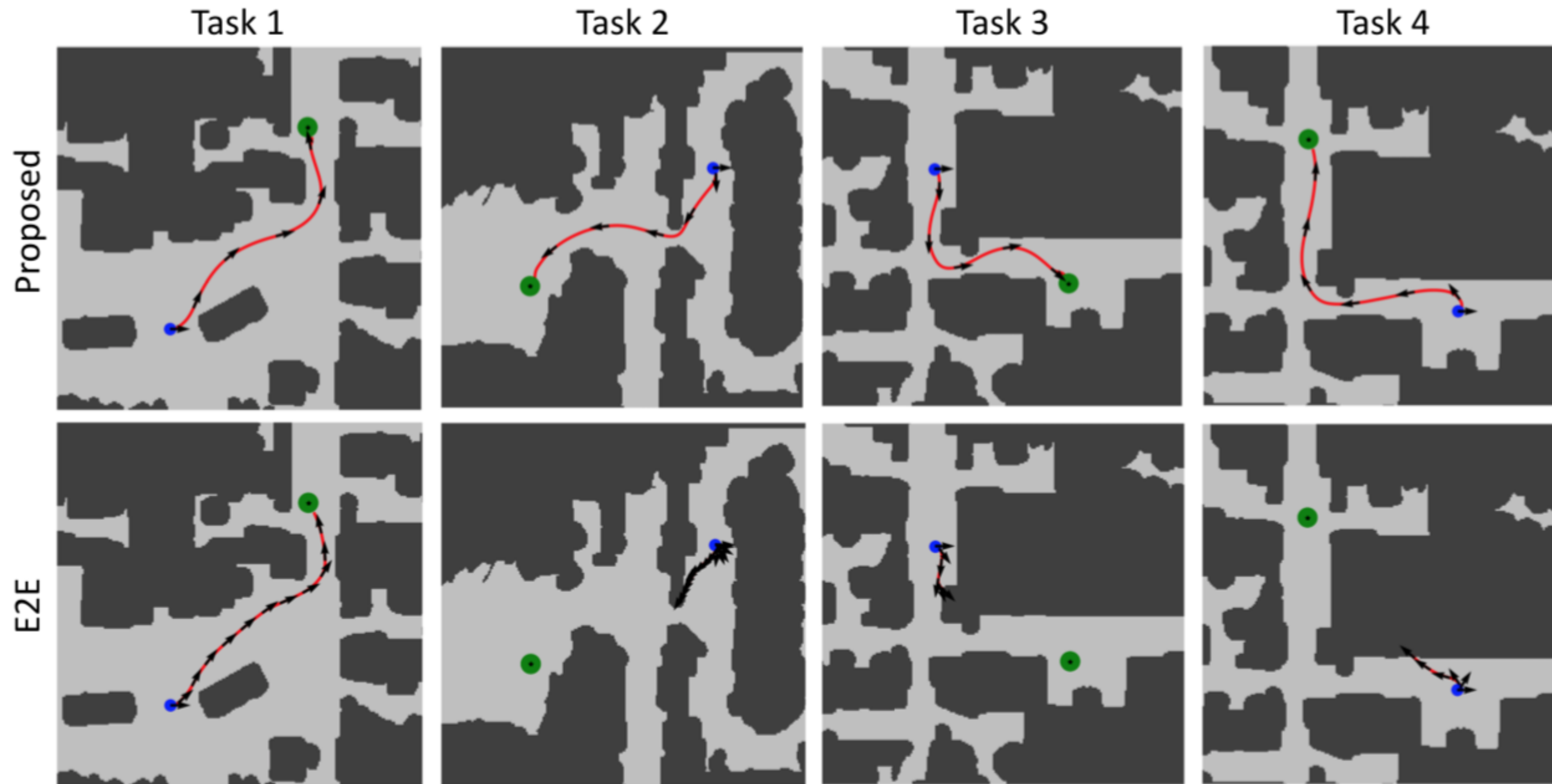


Average jerk along the trajectory ( $m/s^3$ ):

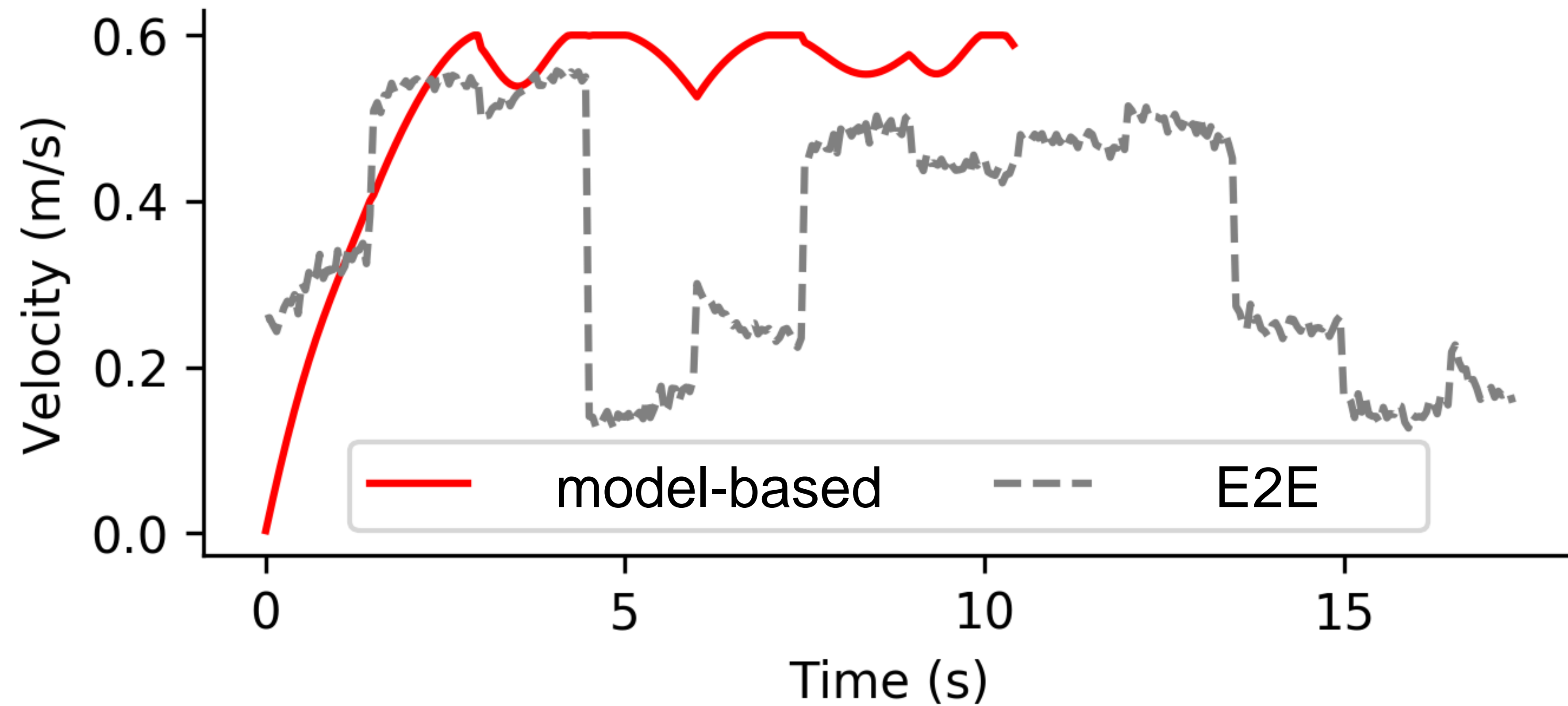


<b>Agent</b>	<b>Input</b>	<b>Success (%)</b>	<b>Time taken (s)</b>	<b>Acceleration (<math>m/s^2</math>)</b>	<b>Jerk (<math>m/s^3</math>)</b>
Expert	Full map	100	10.78 $\pm$ 2.64	0.11 $\pm$ 0.03	0.36 $\pm$ 0.14
LB-WayPtNav (our)	RGB	80.65	11.52 $\pm$ 3.00	0.10 $\pm$ 0.04	0.39 $\pm$ 0.16
End To End	RGB	58.06	19.16 $\pm$ 10.45	0.23 $\pm$ 0.02	8.07 $\pm$ 0.94
Mapping (memoryless)	Depth	86.56	10.96 $\pm$ 2.74	0.11 $\pm$ 0.03	0.36 $\pm$ 0.14
Mapping	Depth + Spatial Memory	97.85	10.95 $\pm$ 2.75	0.11 $\pm$ 0.03	0.36 $\pm$ 0.14

# Success Rate

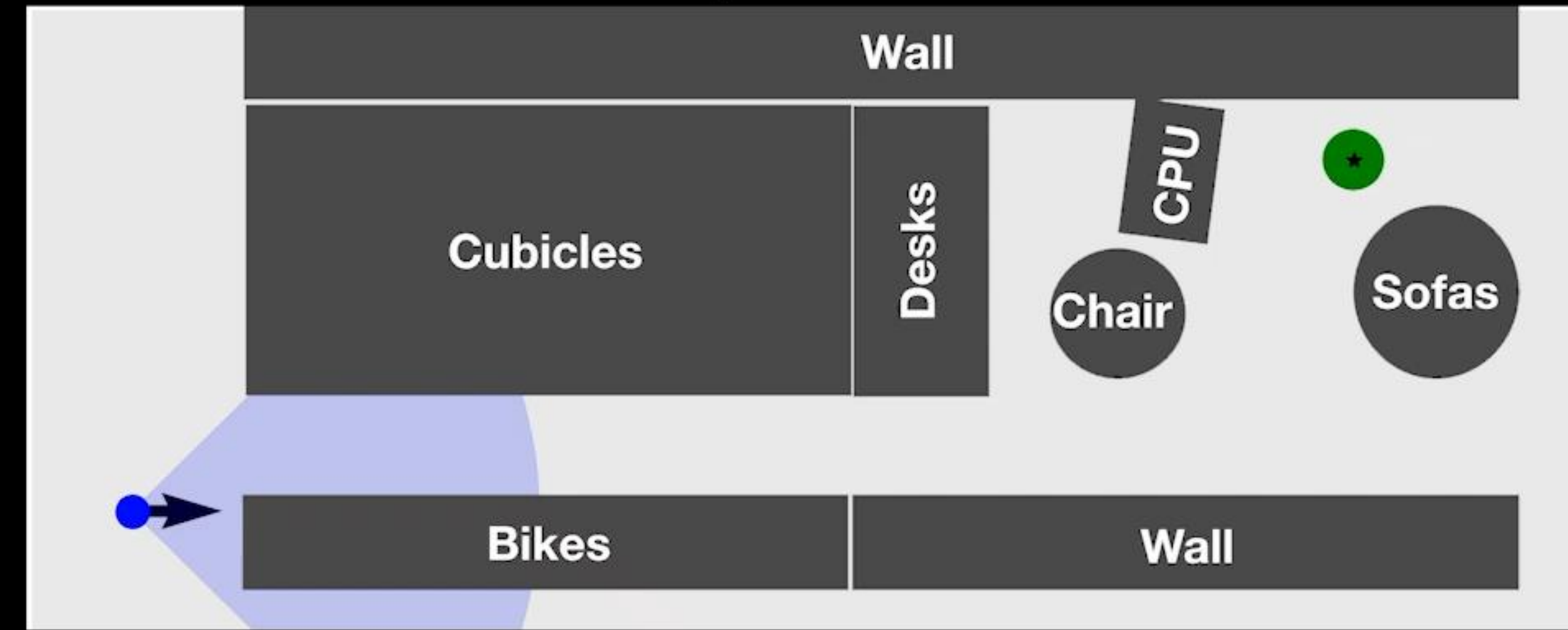


# Control Profile





# Top View



# First Person View



# Third Person View





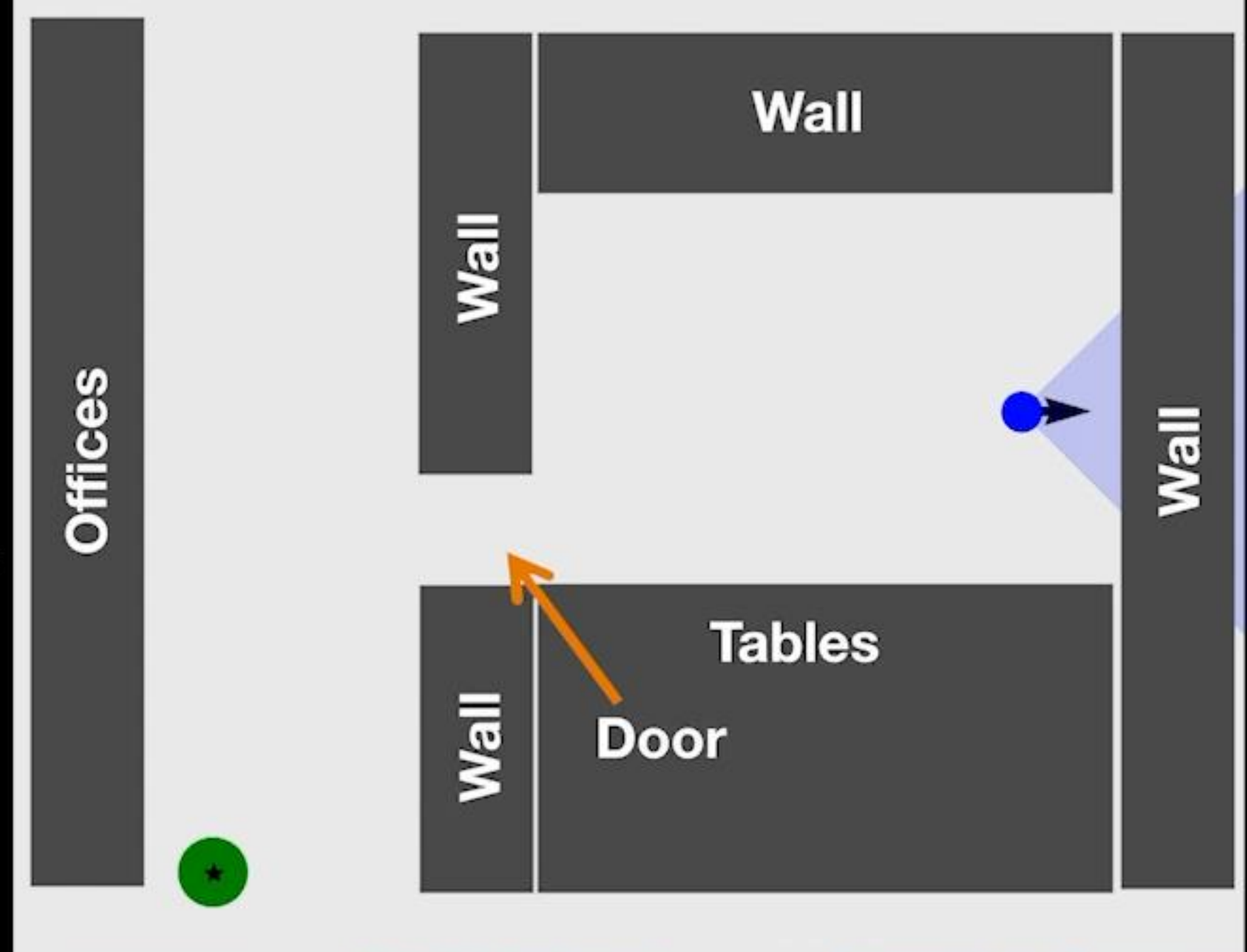
First Person View



Third Person View

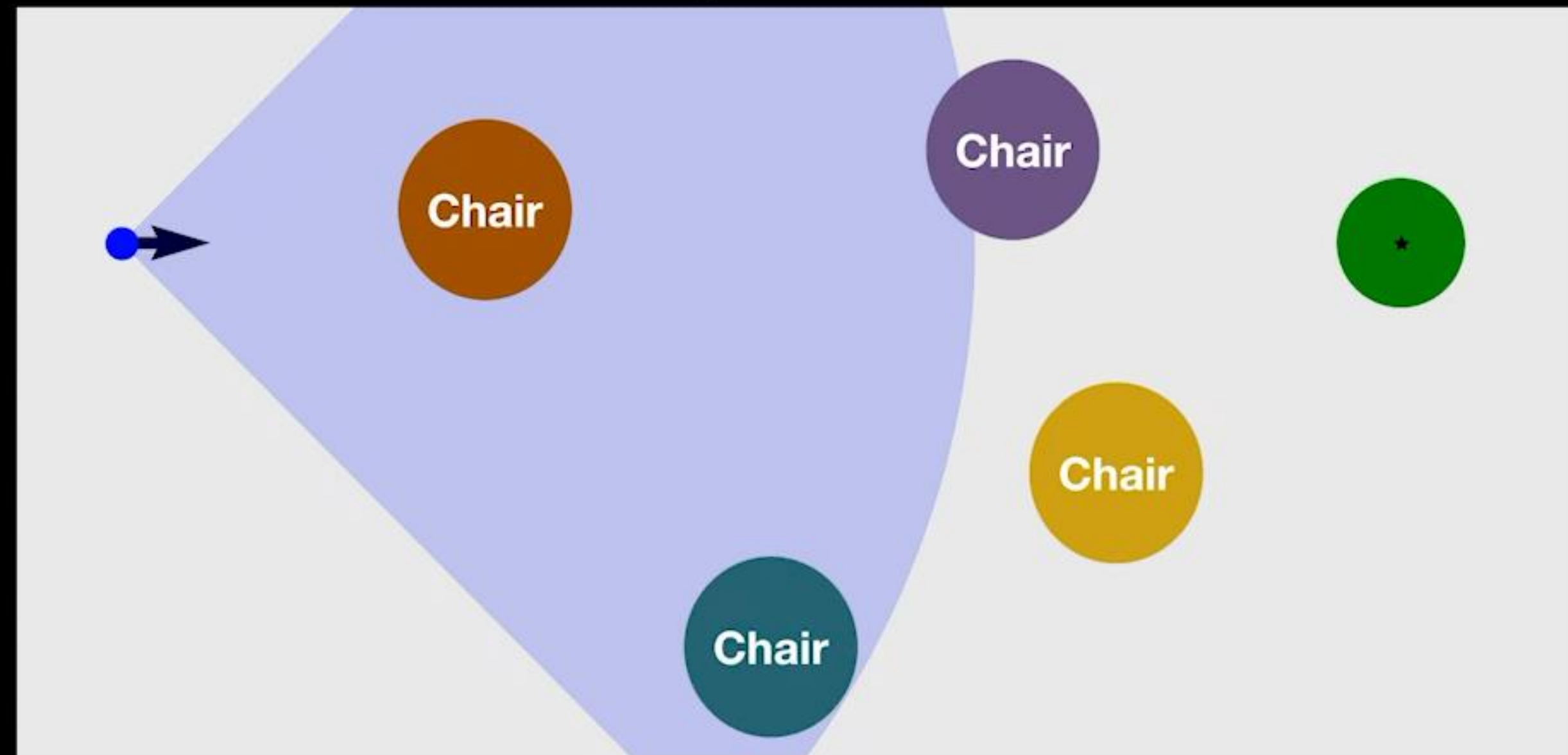


Top View





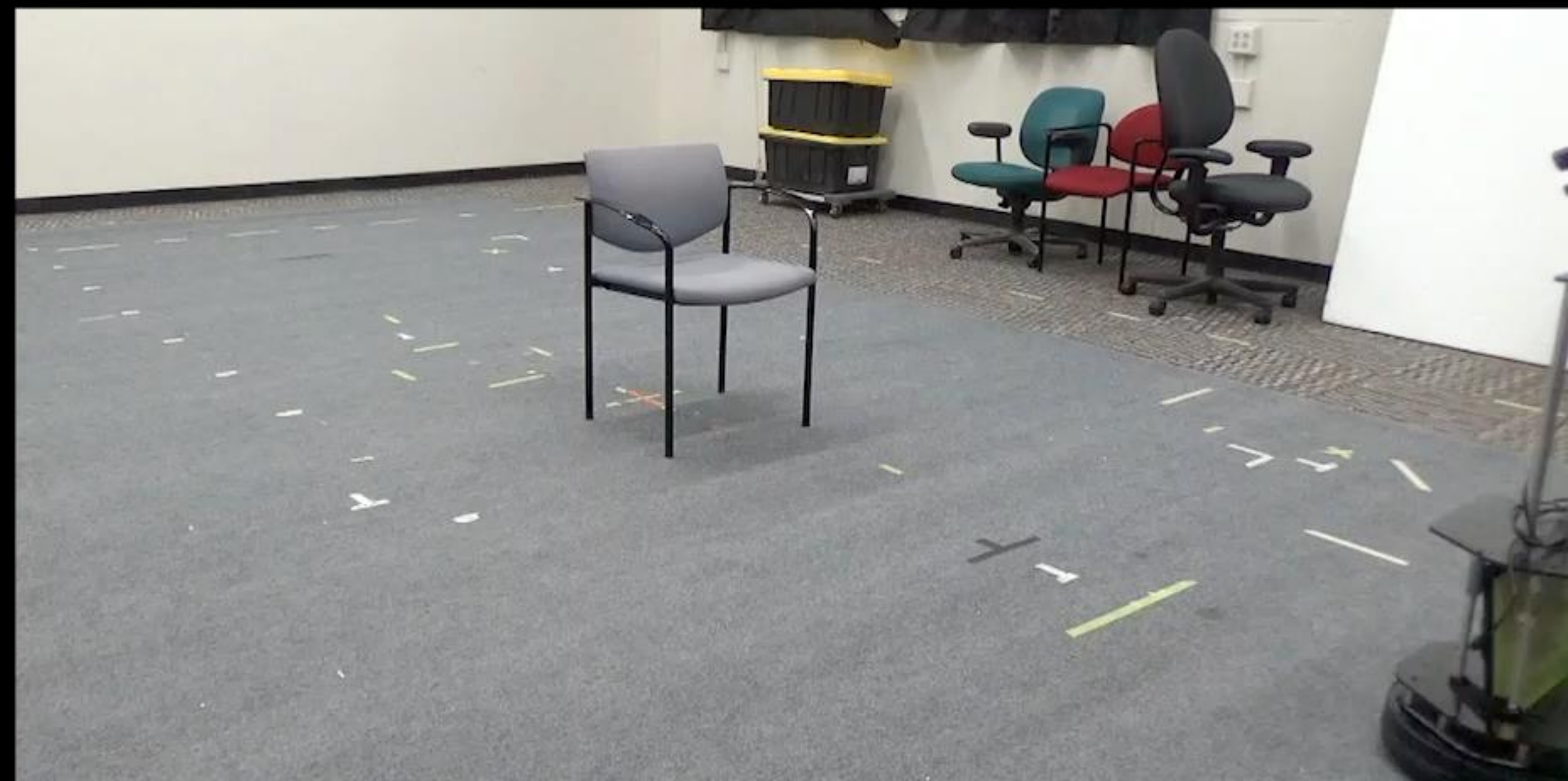
# Top View



# First Person View



# Third Person View





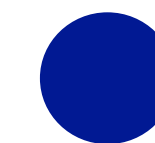
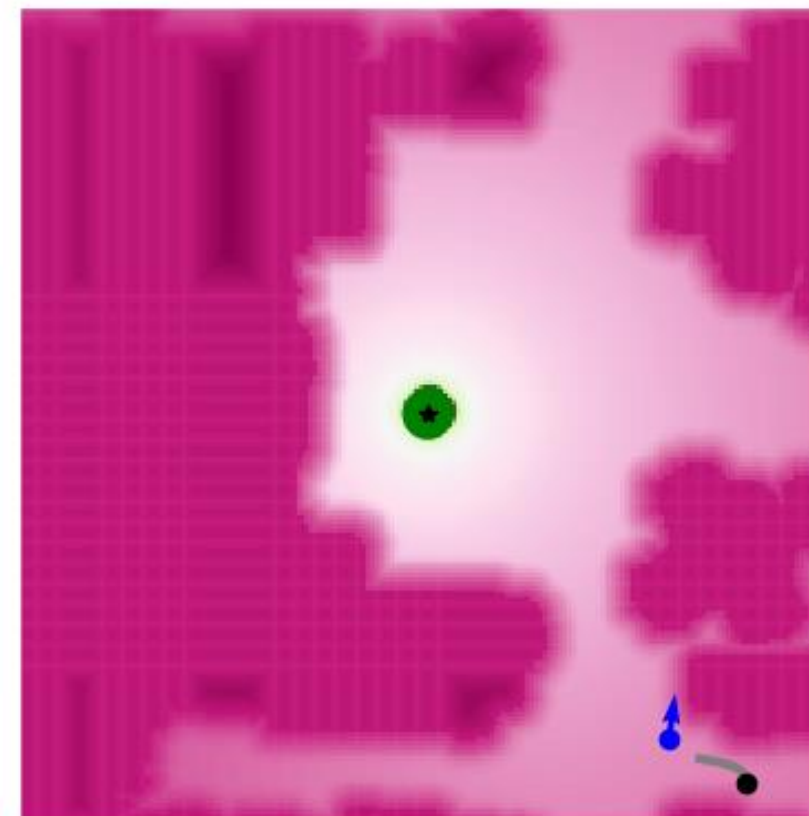
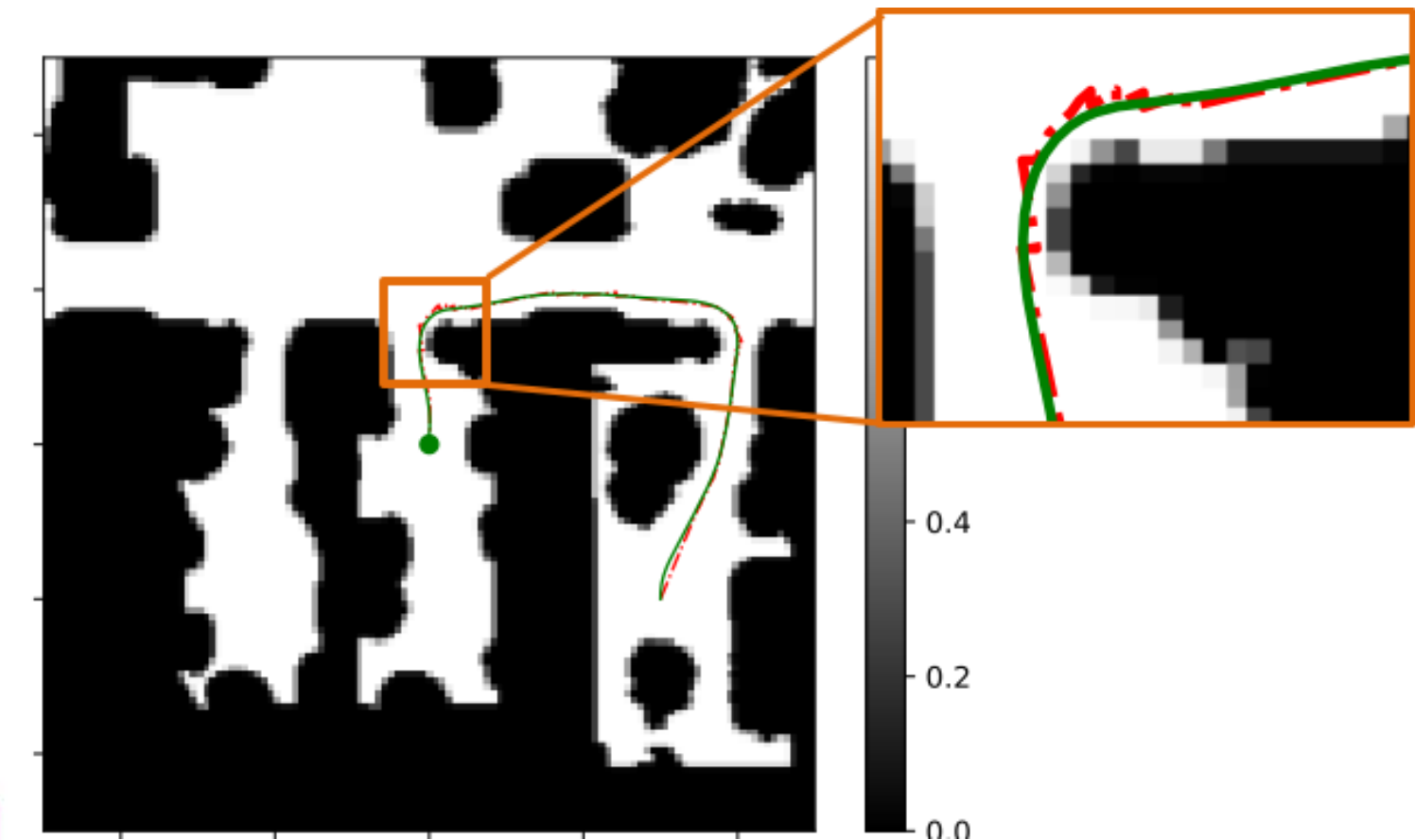
<b>Agent</b>	<b>Input</b>	<b>Success (%)</b>	<b>Time taken (<math>s</math>)</b>	<b>Acceleration (<math>m/s^2</math>)</b>	<b>Jerk (<math>m/s^3</math>)</b>
LB-WayPtNav (our)	RGB	95	$22.93 \pm 2.38$	$0.09 \pm 0.01$	$3.01 \pm 0.38$
End To End	RGB	50	$33.88 \pm 3.01$	$0.19 \pm 0.01$	$6.12 \pm 0.18$
Mapping (memoryless)	RGB-D	0	N/A	N/A	N/A
Mapping	RGB-D + Spatial Memory	40	$22.13 \pm 0.54$	$0.11 \pm 0.01$	$3.44 \pm 0.21$

# Some lessons learned

- Data representation is important
- Optimal control can be too optimal
- Waypoint representation



vs



vs





# More lessons learned

- Building on existing NN architectures
- Image and perspective distortions during training
- RL on supervised learning

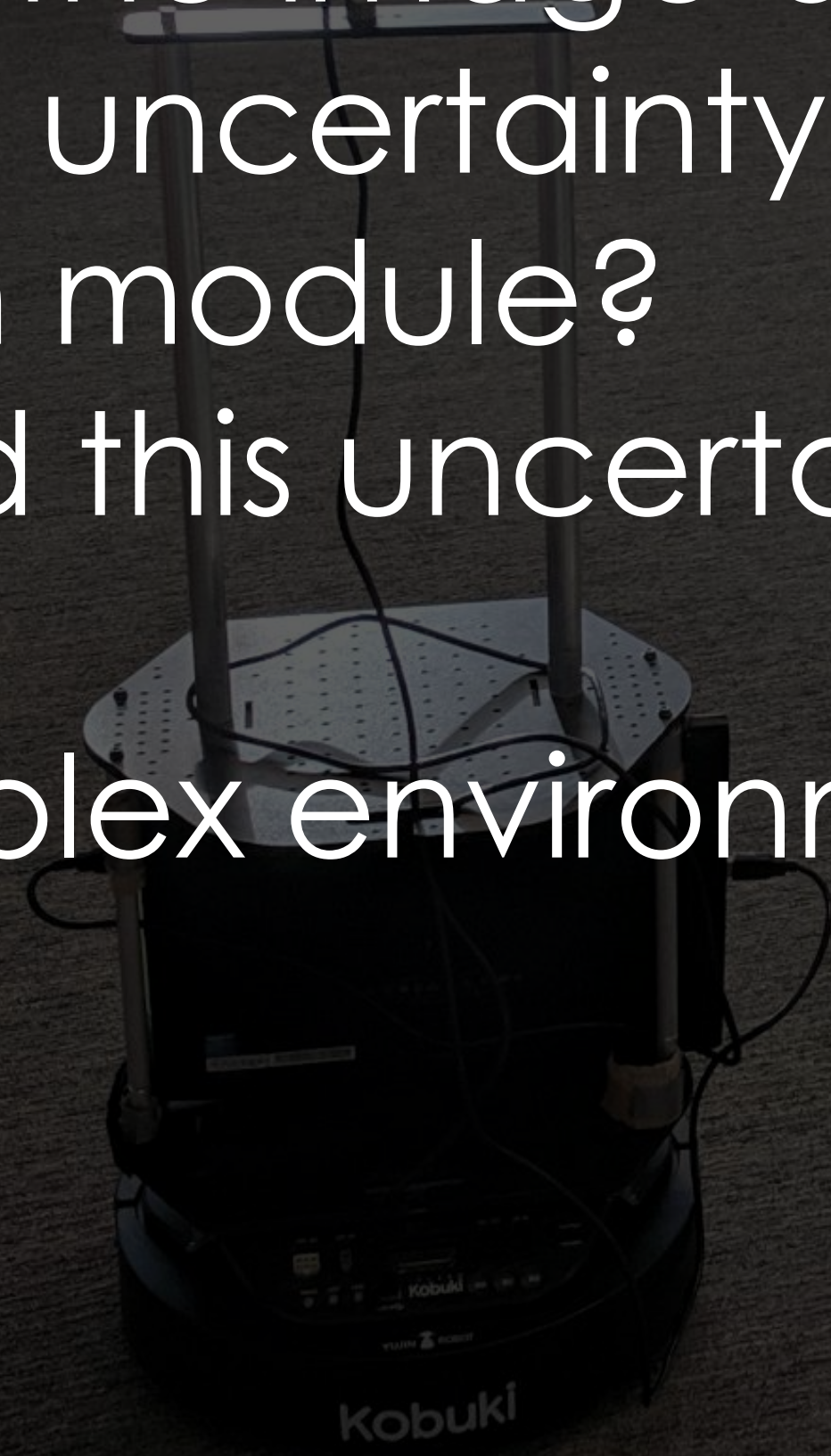






# Safety Challenges

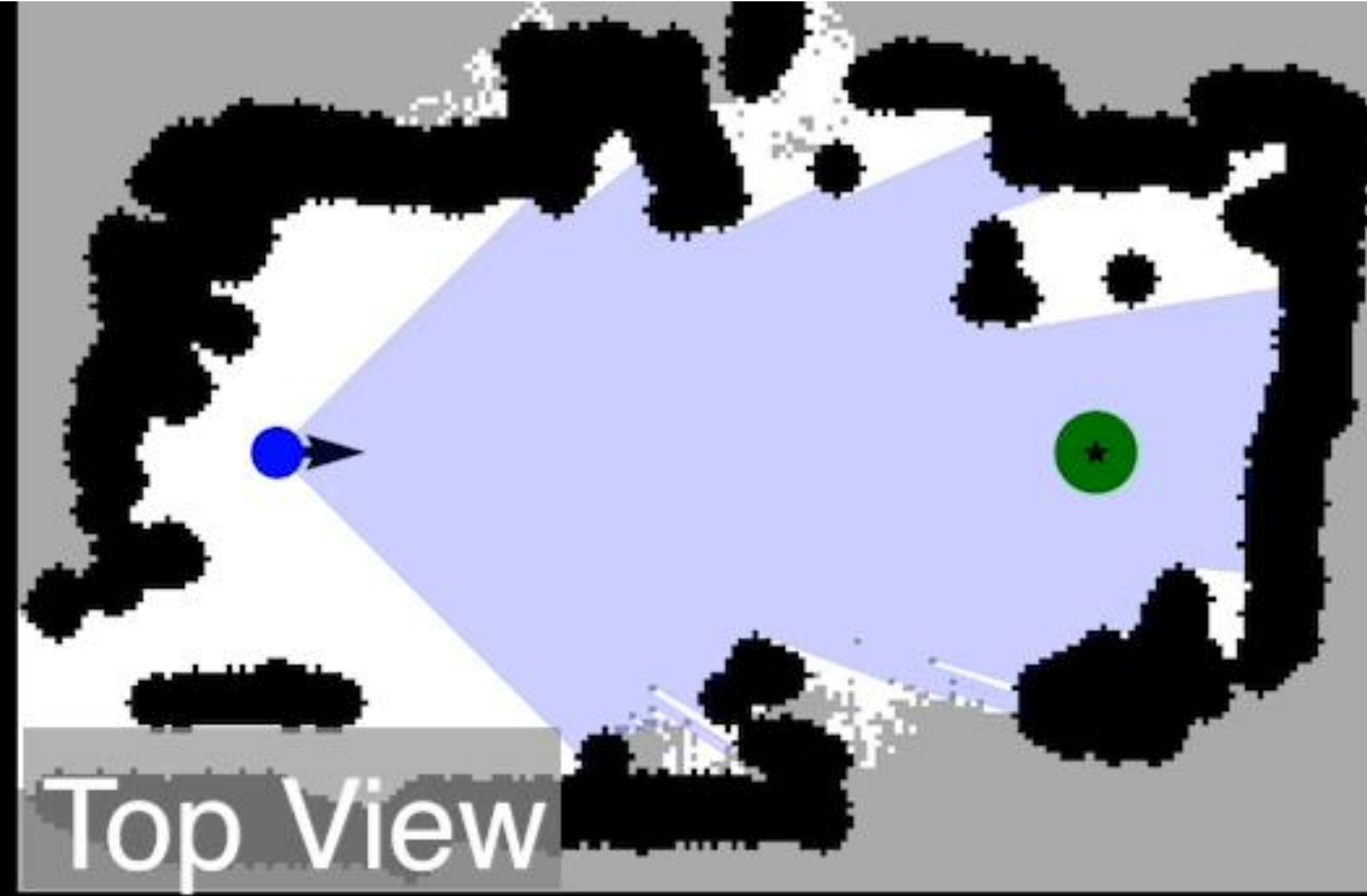
- Monitor: Is the image data in the training distribution?
- What is the uncertainty around the output of the perception module?
- How should this uncertainty affect the planning and control?
- More complex environments?





# LB-WayPtNav

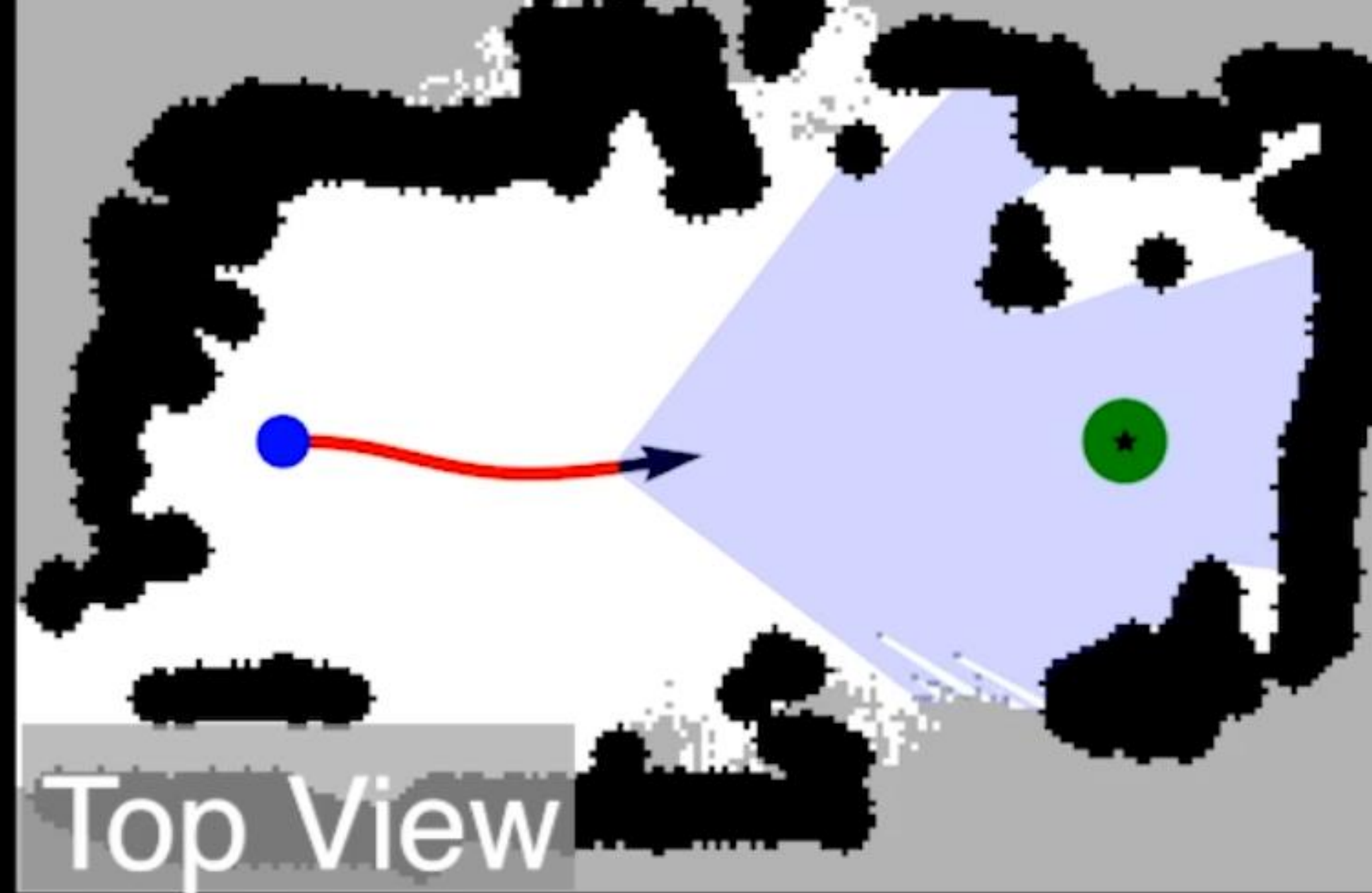
## Experiment 1 (1x)





# LB-WayPtNav

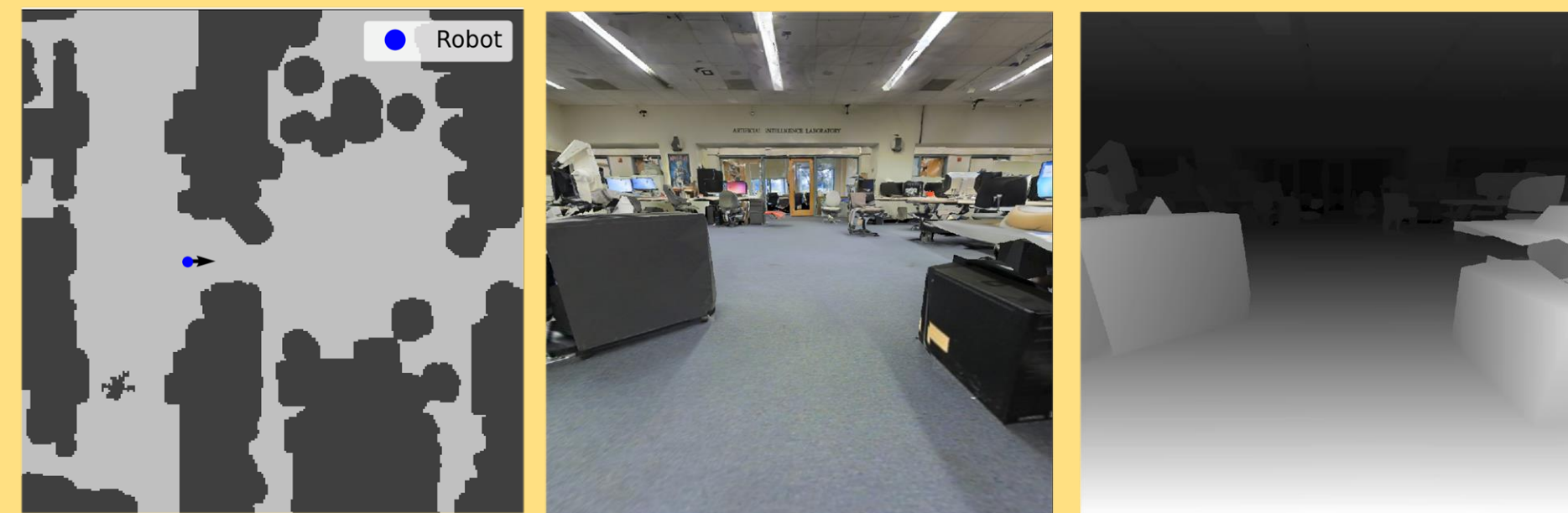
## Experiment 1 (1x)



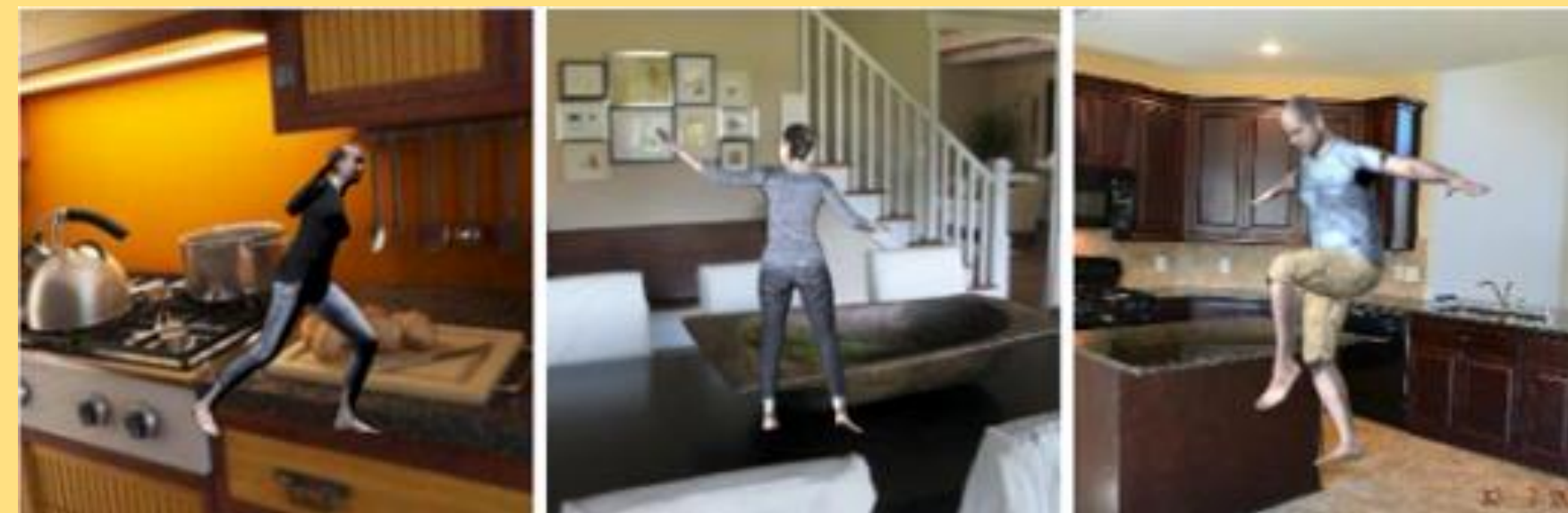


# HumANav

SD3DIS



SURREAL





# HumANav

Robot State

$[x, y, \theta]$

Human State & Control

$[x, y, \theta, v, \omega]$

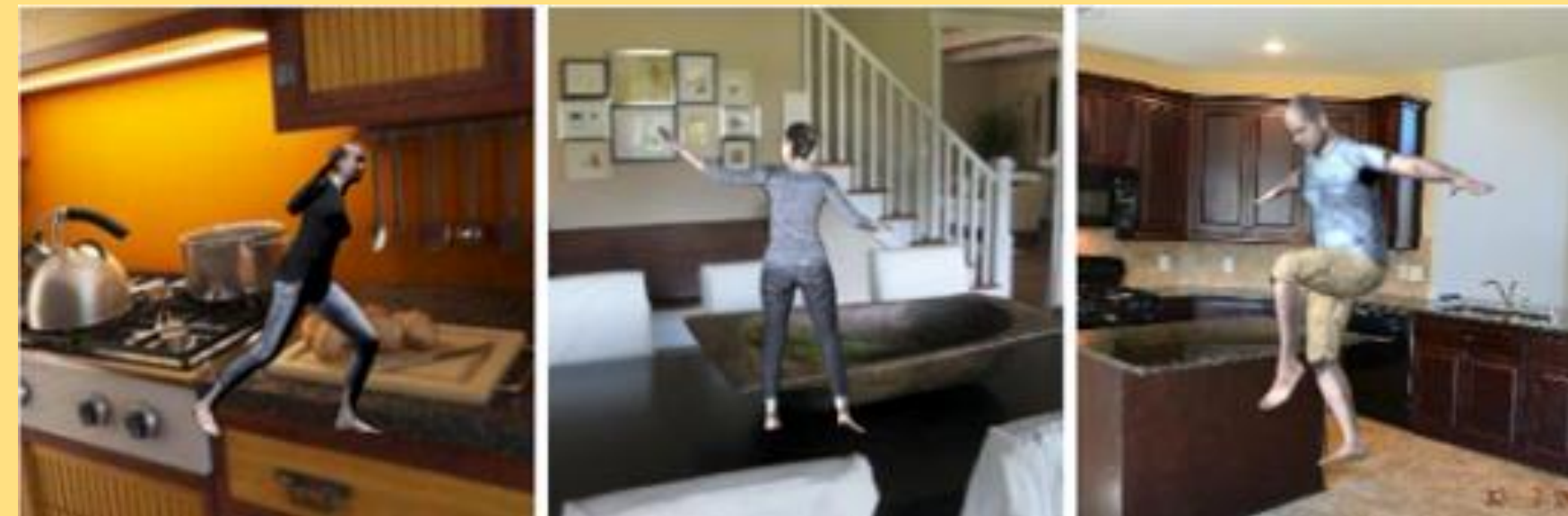
Human Identity

- Gender
- Texture (Clothing, Skin Color, Facial Features)
- Body Shape (Tall, Short, etc.)

SD3DIS



SURREAL





# HumANav

Robot State

$[x, y, \theta]$

Human State & Control

$[x, y, \theta, v, \omega]$

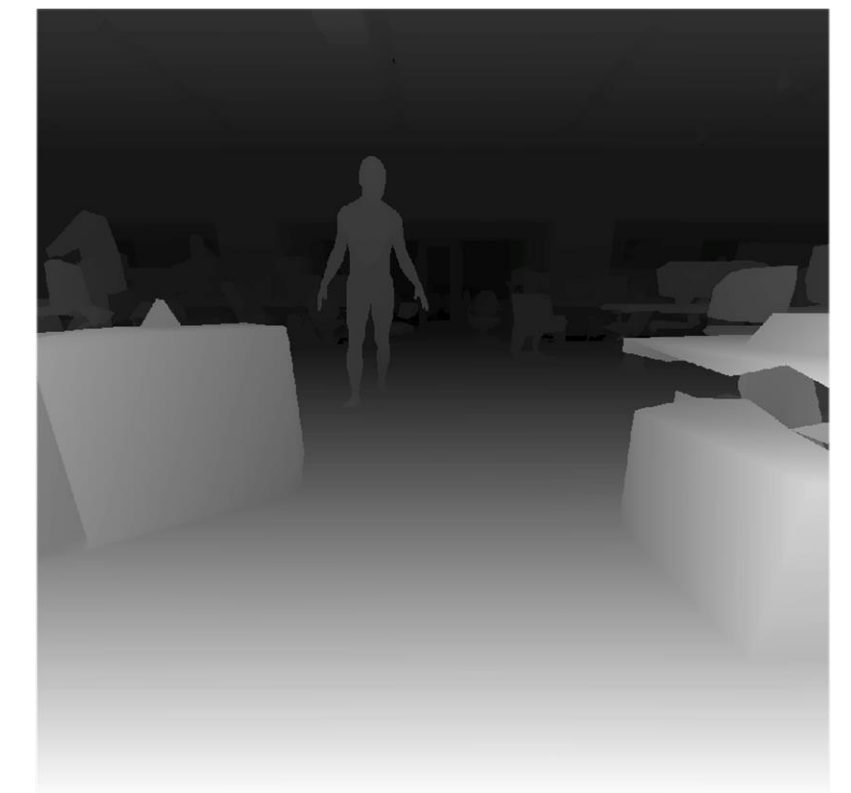
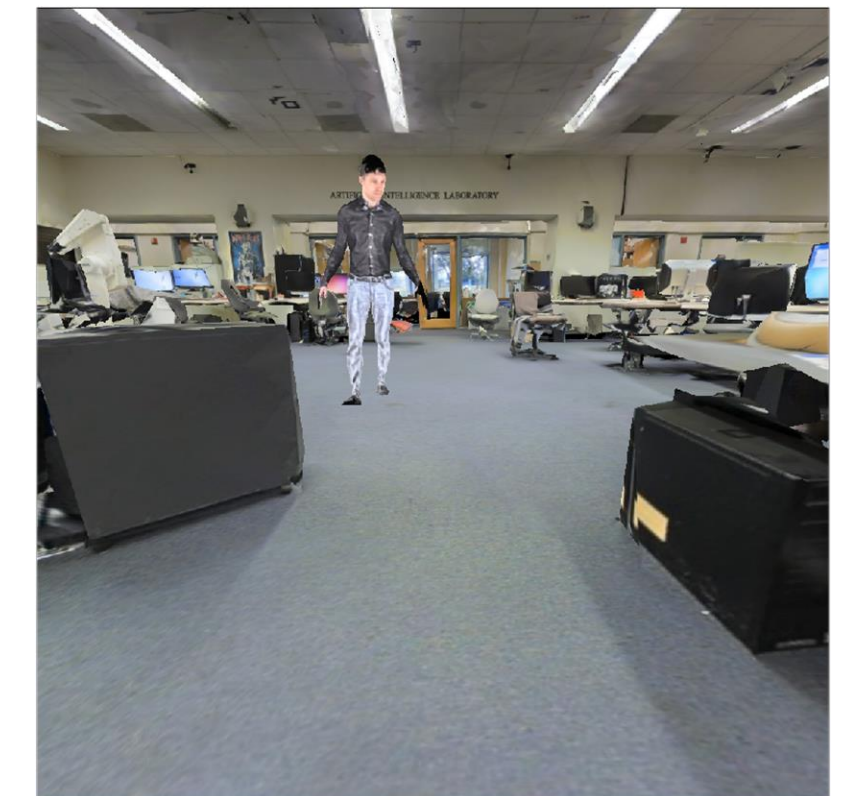
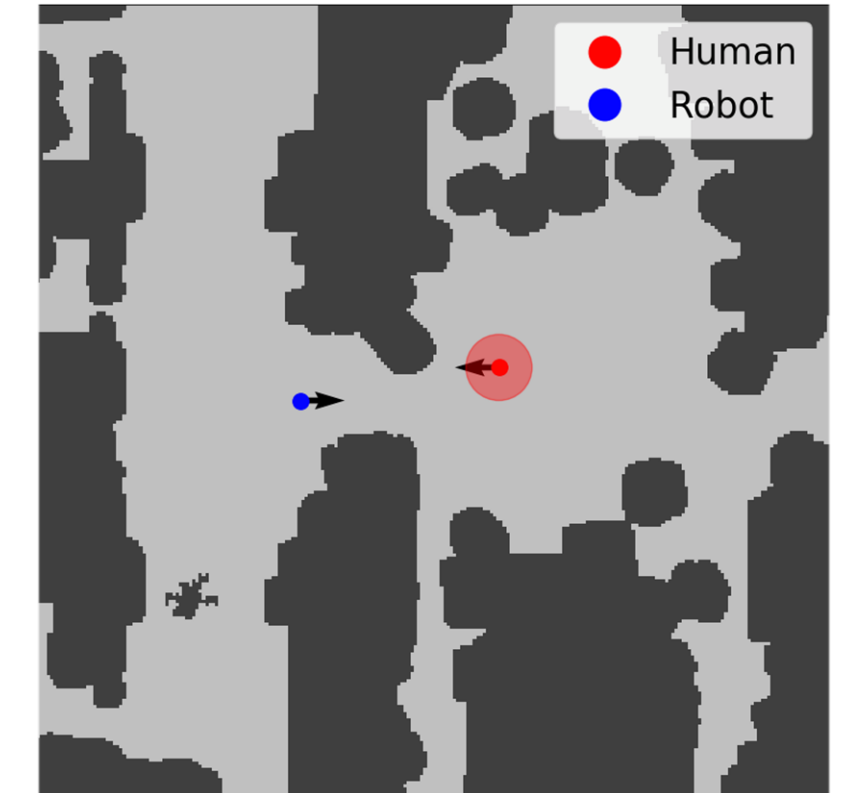
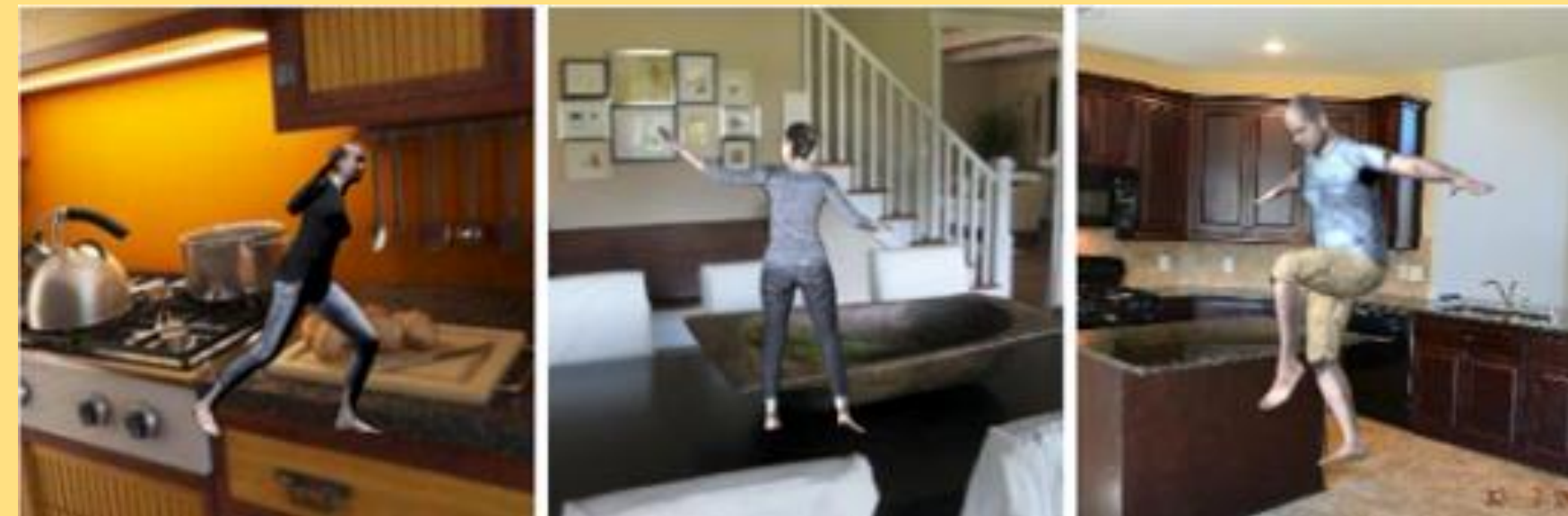
Human Identity

- Gender
- Texture (Clothing, Skin Color, Facial Features)
- Body Shape (Tall, Short, etc.)

SD3DIS



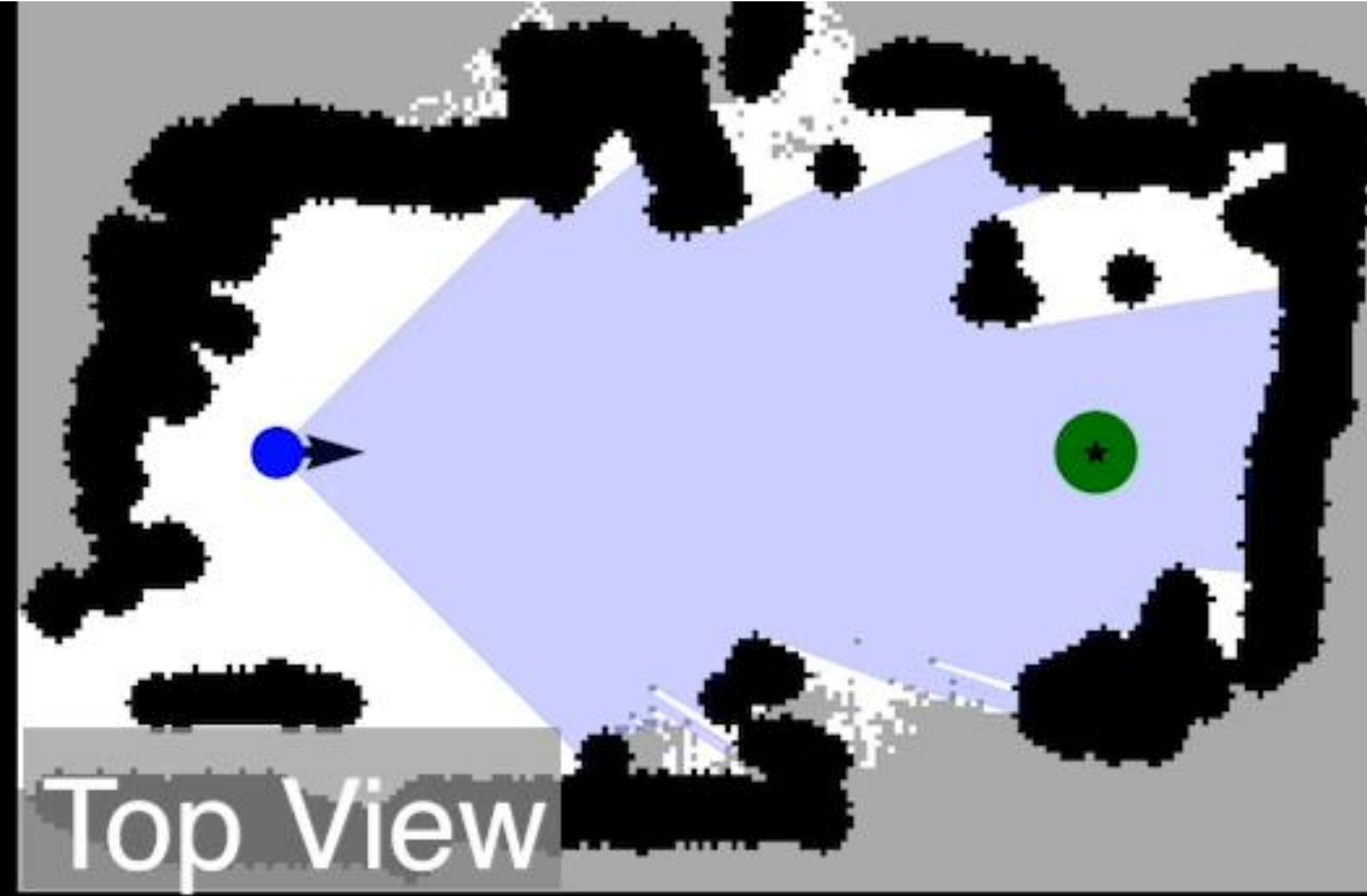
SURREAL





# LB-WayPtNav

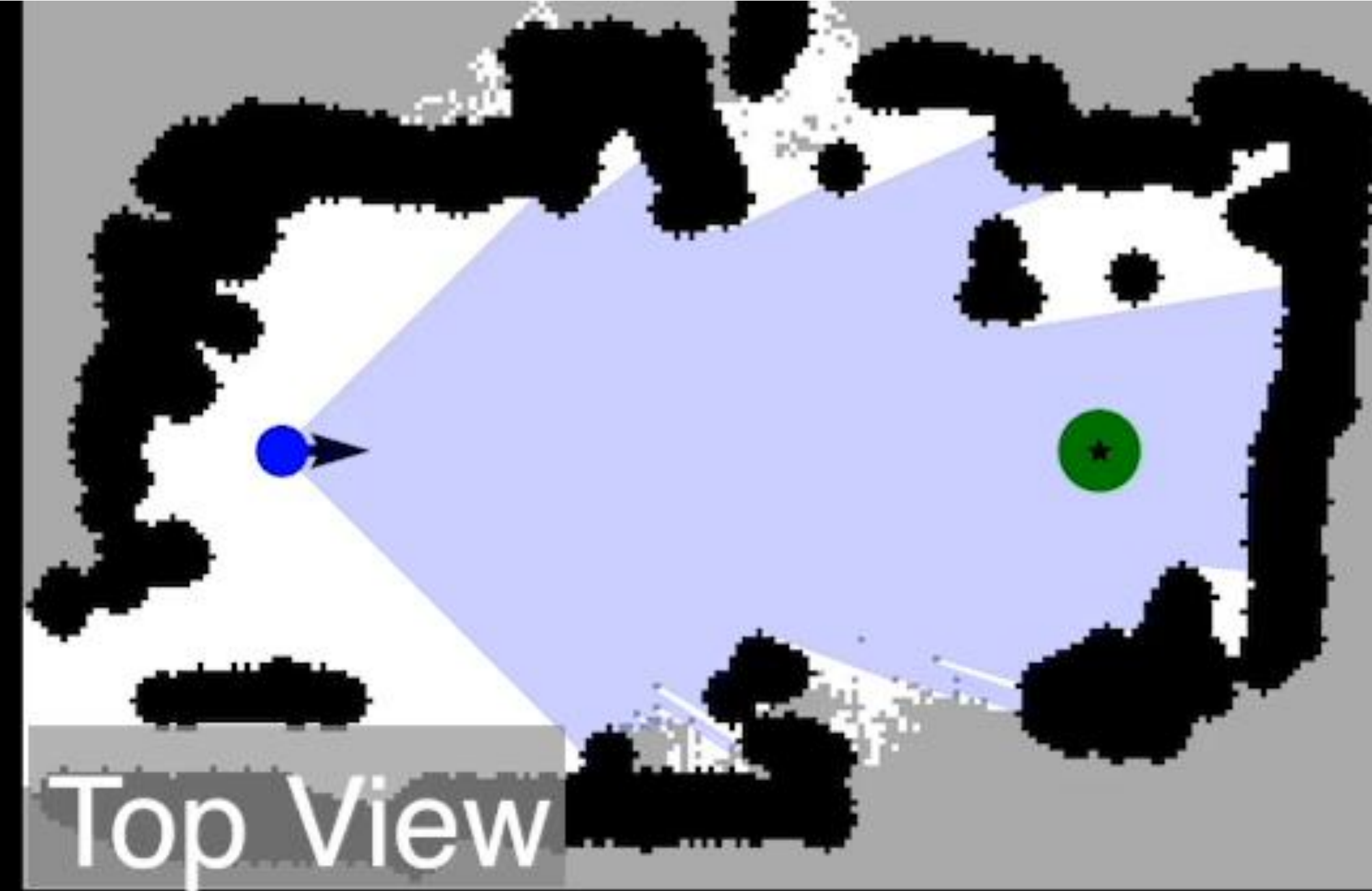
## Experiment 1 (1x)





# LB-WayPtNav-DH

## Experiment 1 (1x)





# Summary

Incorporating perception in the control loop

Supervising learning using optimal control

- a perception-planning-control pipeline
- comparison with a more traditional SLAM pipeline
- applied to a vision-based navigation task

Models of human motion

Challenges for safe control