

CS 4803 / 7643: Deep Learning

Topics: RL and Robotics

- Embodied AI
- Proximal Policy Optimization (PPO)
- Application: Robotics
 - PointGoal Navigation

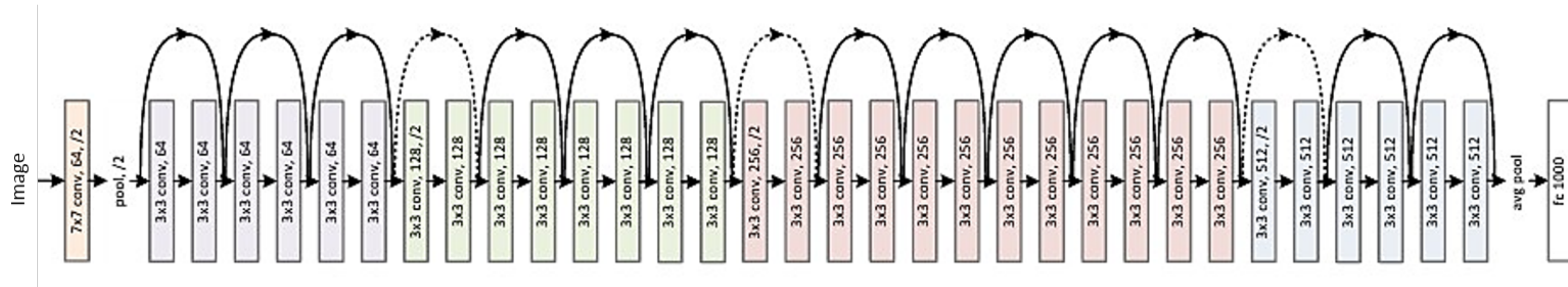
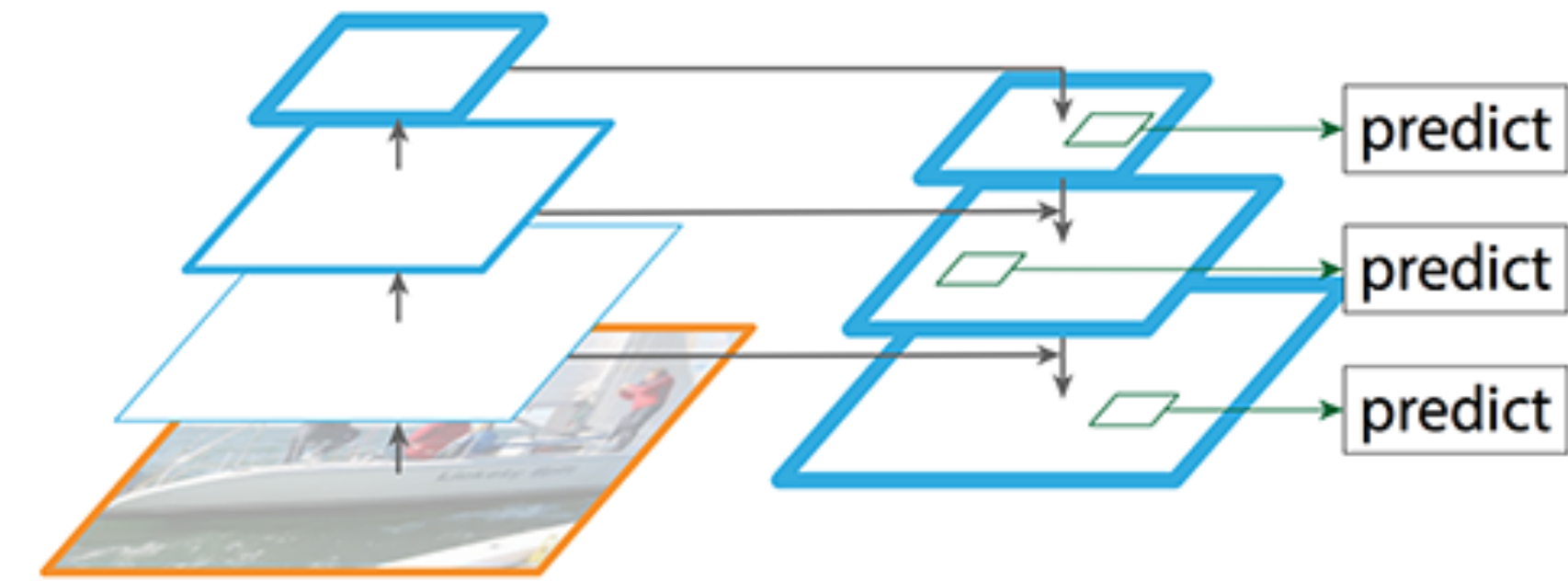
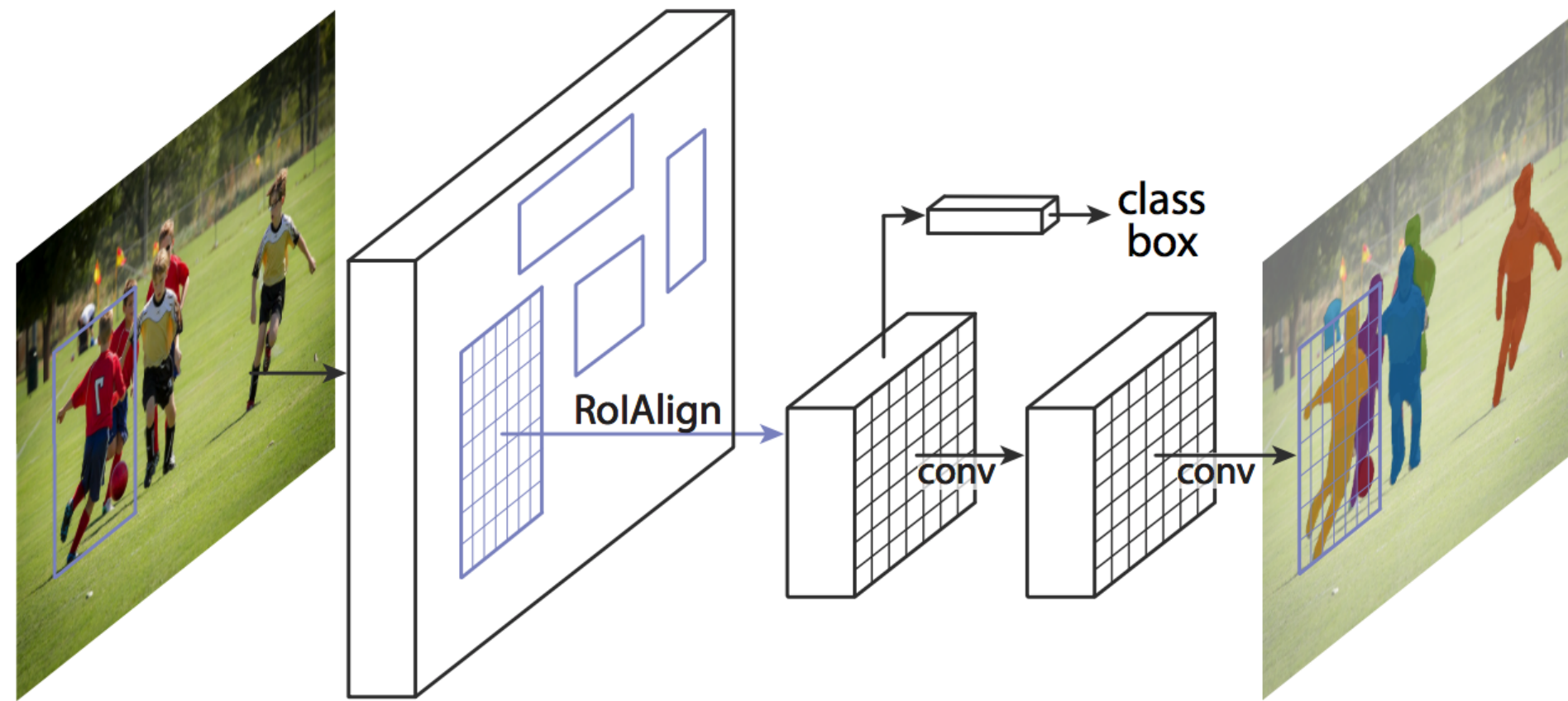
Joanne Truong
Georgia Tech

Lecture Plan

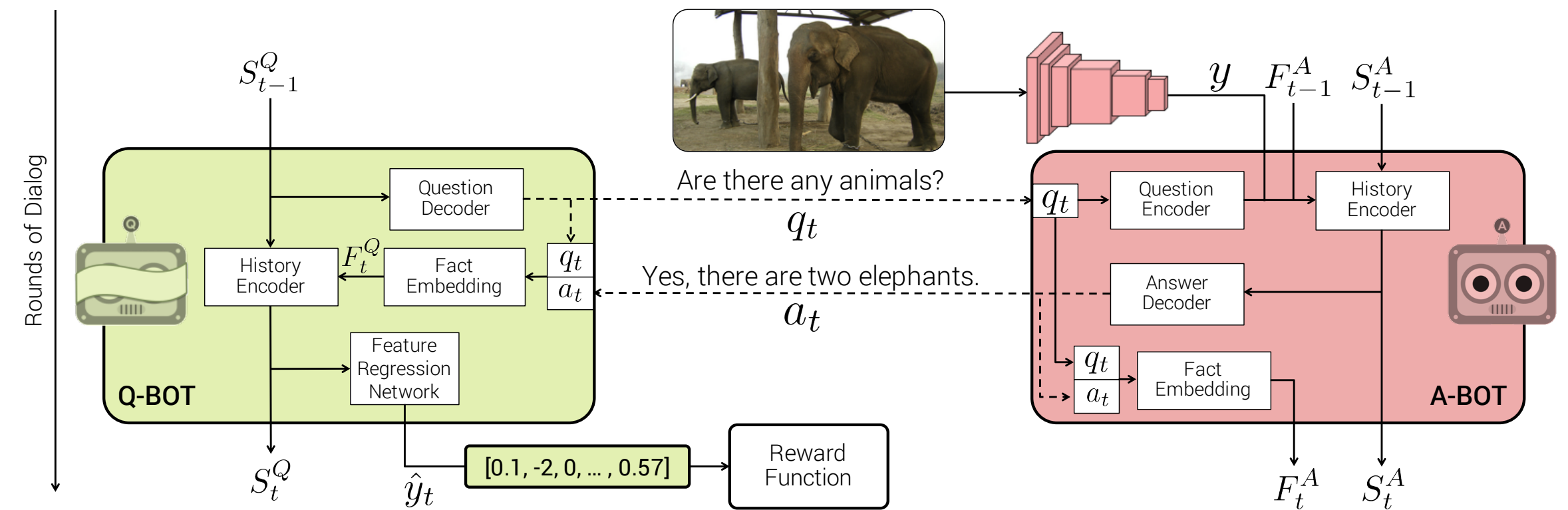
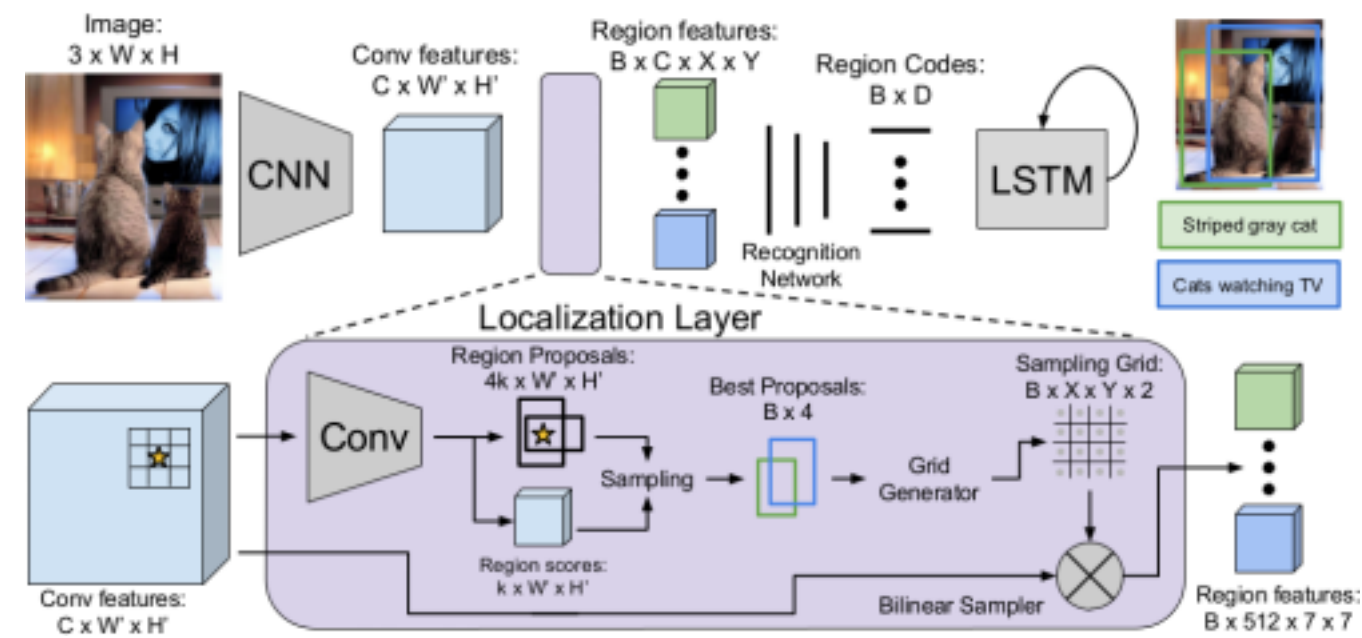
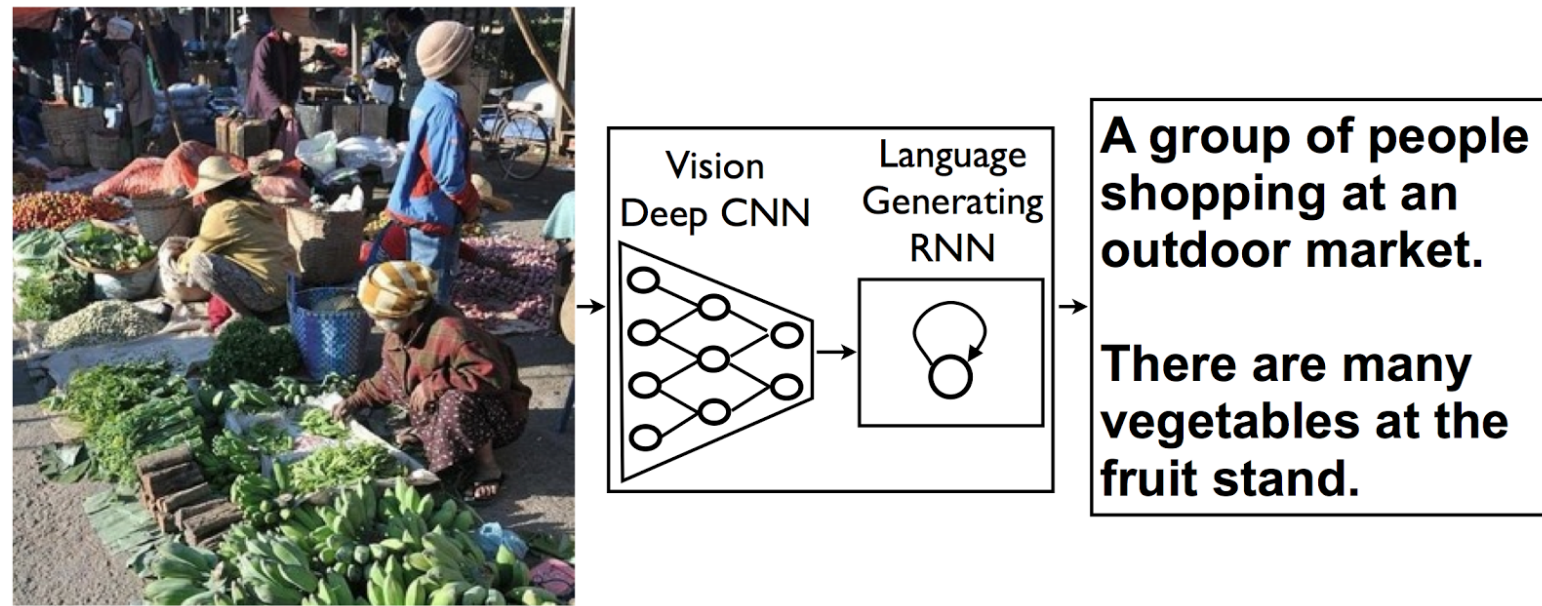
- Embodied AI
- Introduce more advanced RL – Proximal Policy Optimization (PPO)
- Application: Robotics
 - PointGoal Navigation: Combine CNNs, RNNs (LSTMs), and RL together

State-of-the-Art Visual Recognition

State-of-the-Art Visual Recognition



State-of-the-Art Visual Recognition



State-of-the-Art Visual Recognition

Applications

Applications



Applications

Physical agent



Applications

Physical agent
capable of taking
actions in the world



Applications

Physical agent
capable of taking
actions in the world

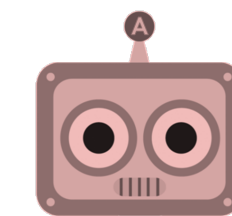


Applications

Physical agent capable of taking actions in the world and talking to humans in natural language

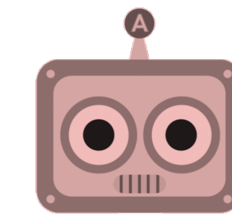


Is there smoke in any room around you?



Yes, in one room

Go there and look for people

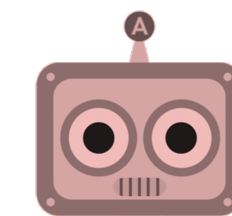


...

Applications

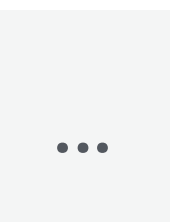
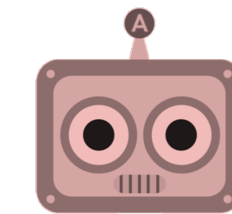


Is there smoke in any room around you?



Yes, in one room

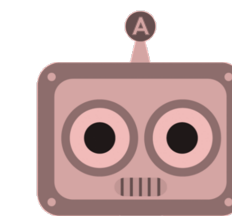
Go there and look for people



Applications

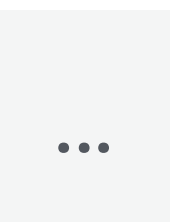
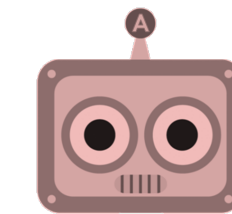


Is there smoke in any room around you?

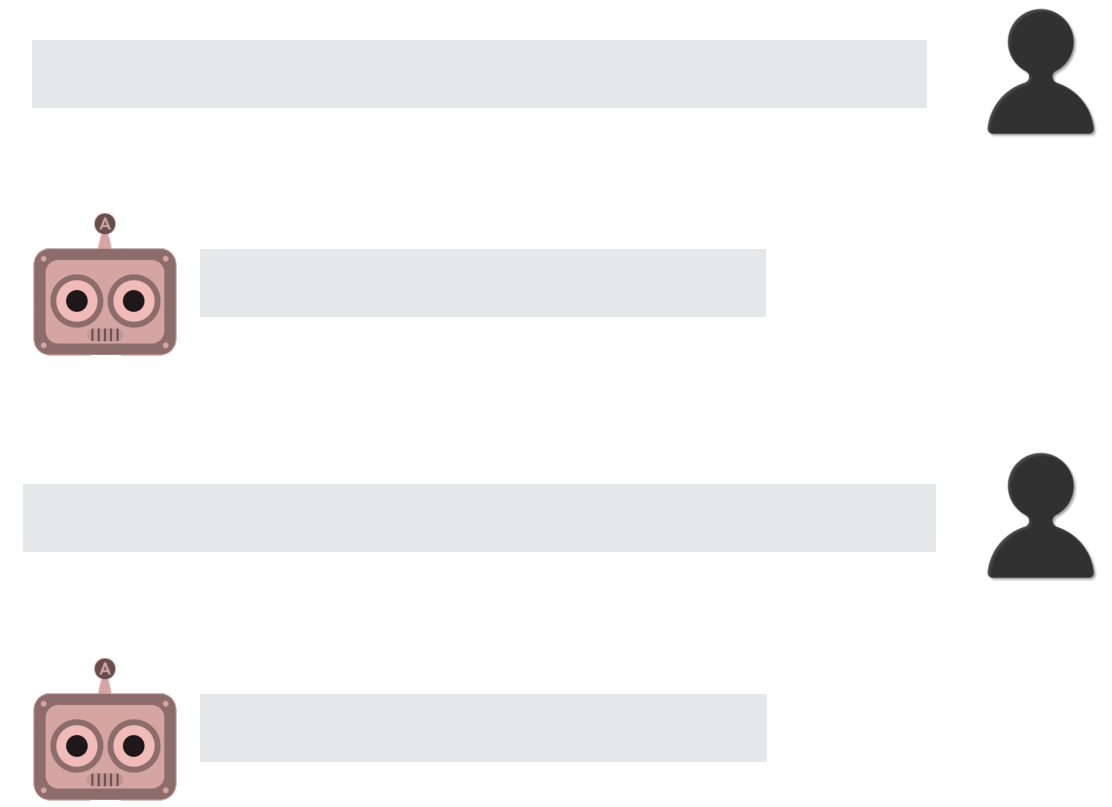


Yes, in one room

Go there and look for people

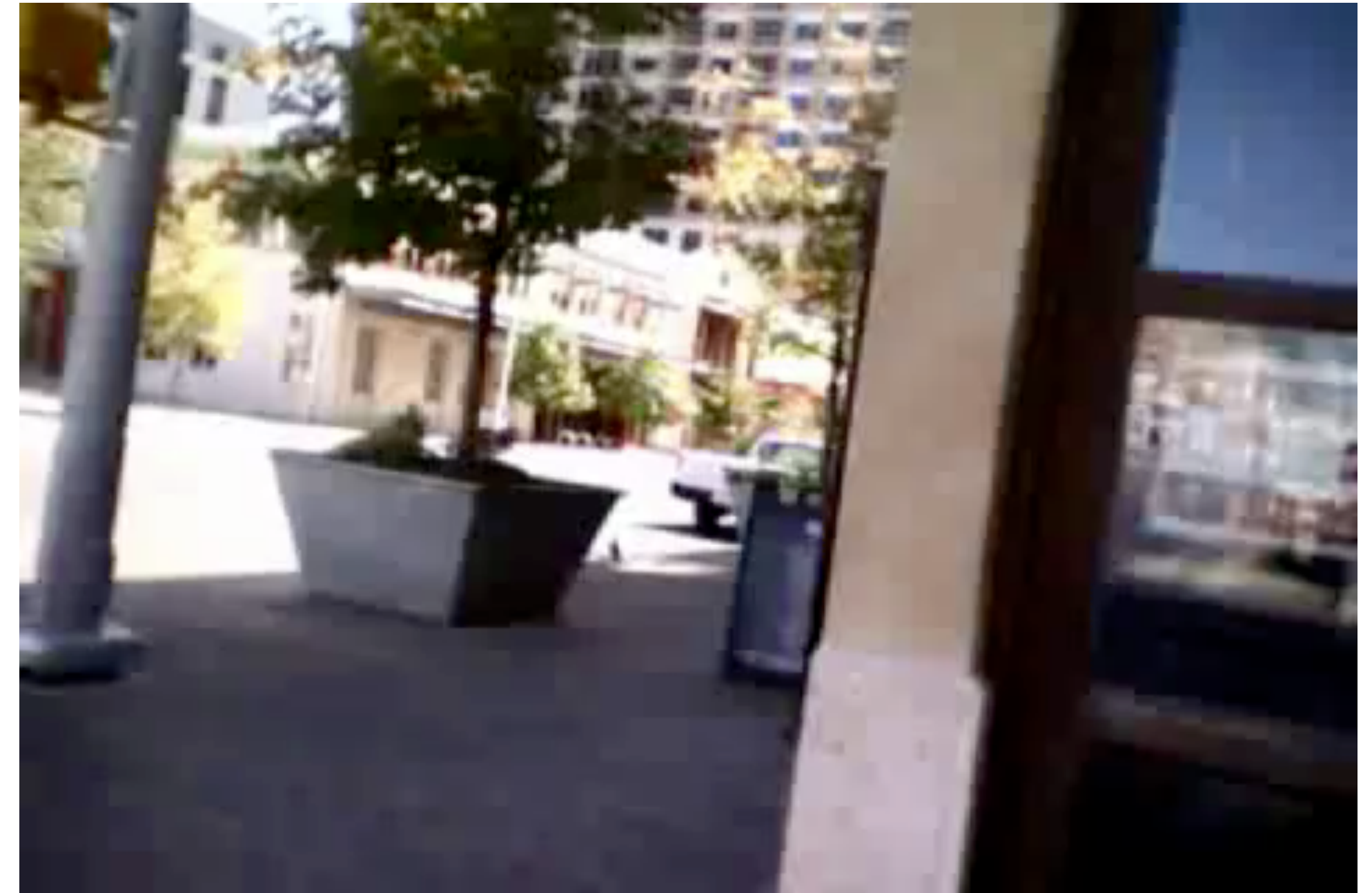


Challenges



Challenges

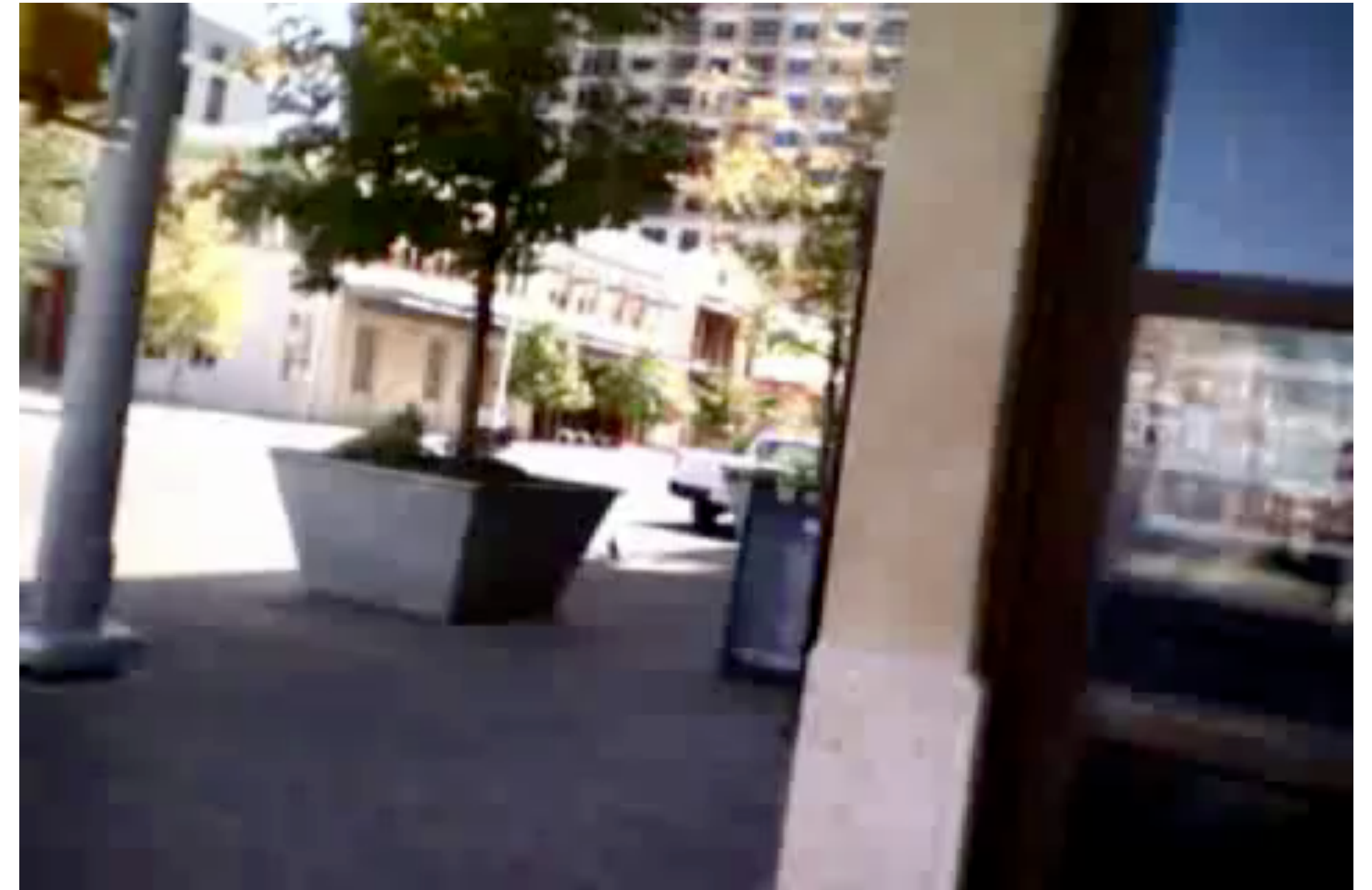
Egocentric vision



No access to well-composed, curated images

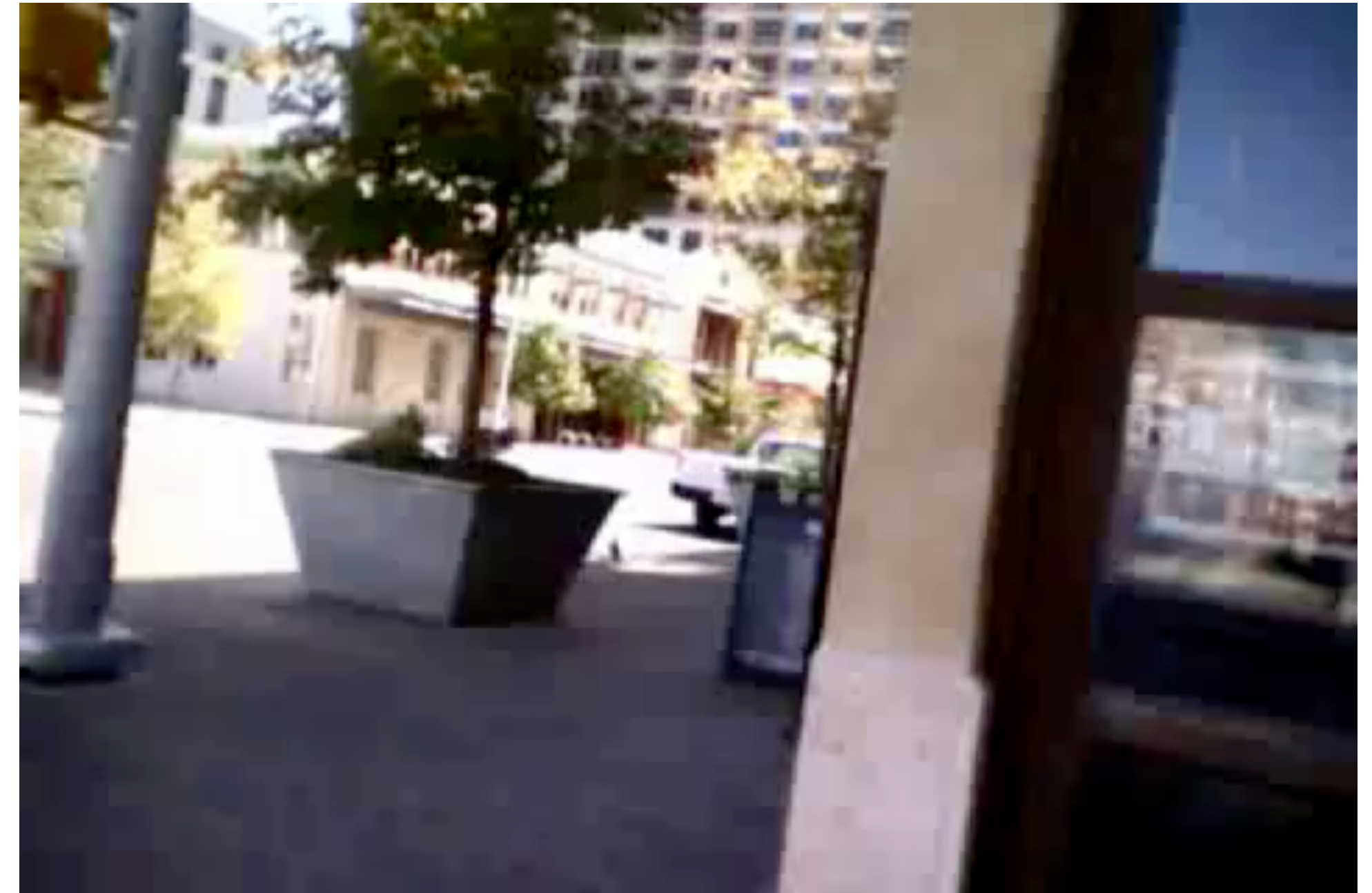
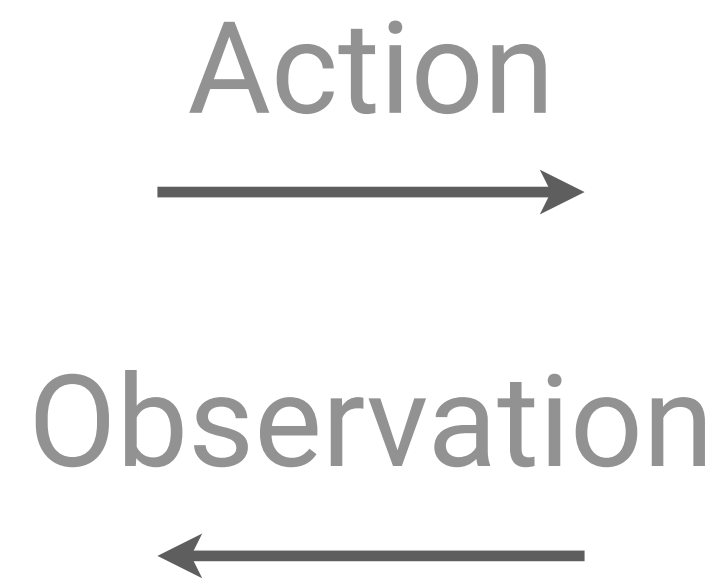
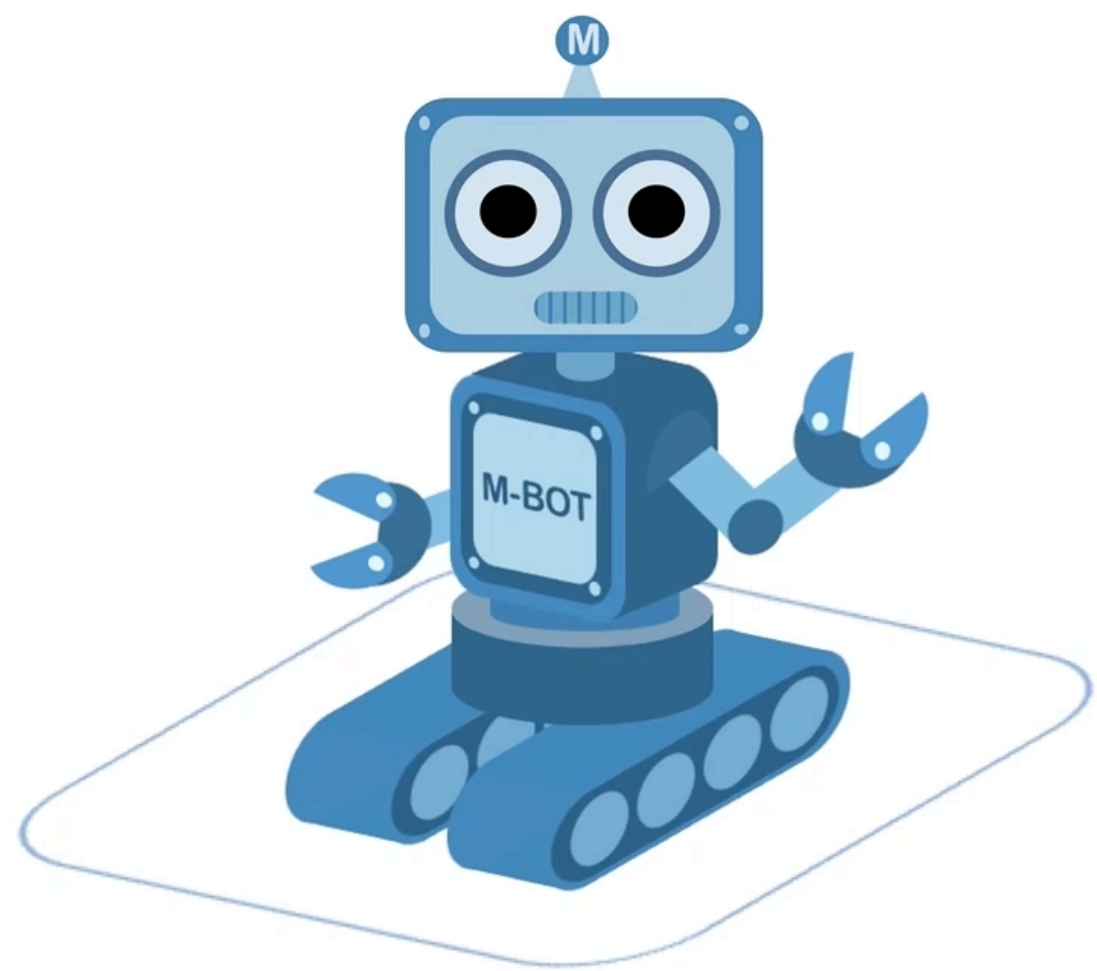
Challenges

Egocentric vision



Challenges

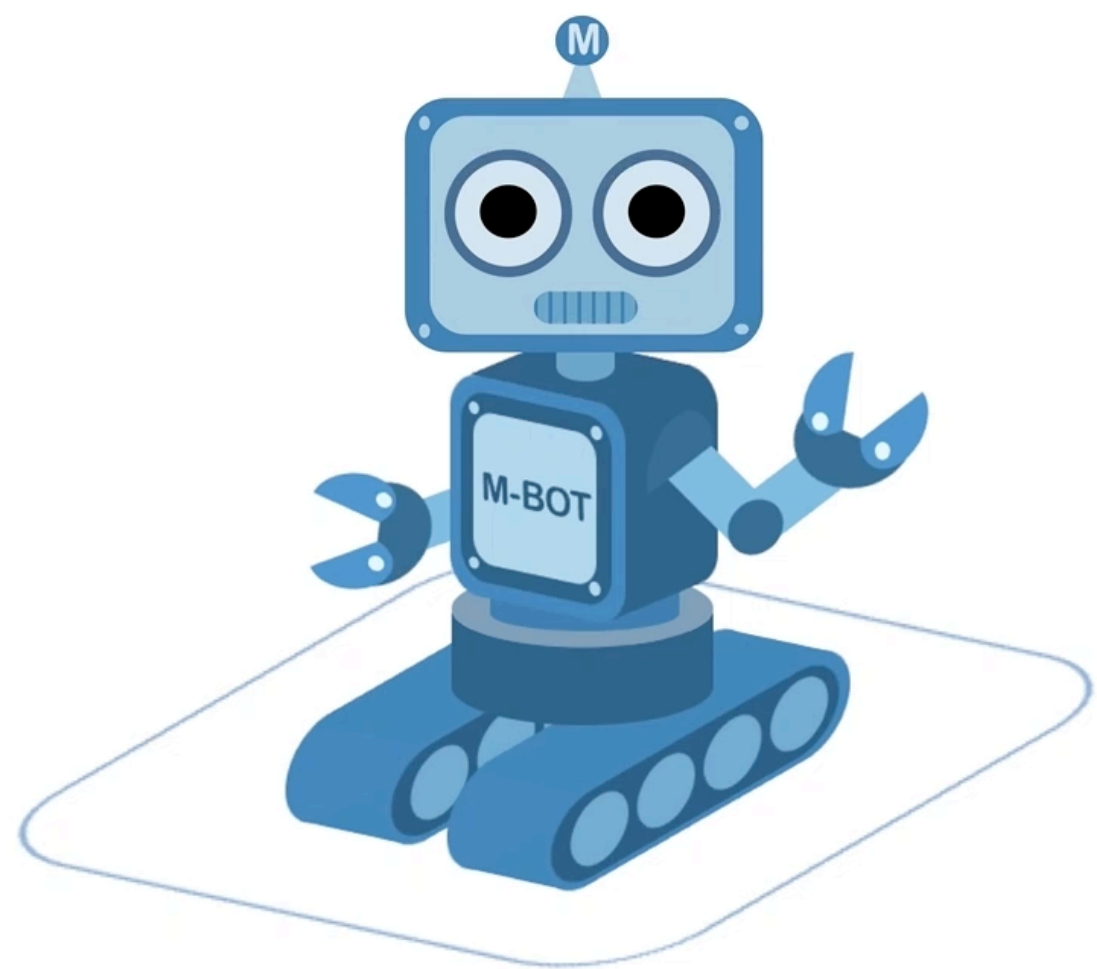
Egocentric vision
Active perception



Agent controls incoming data distribution

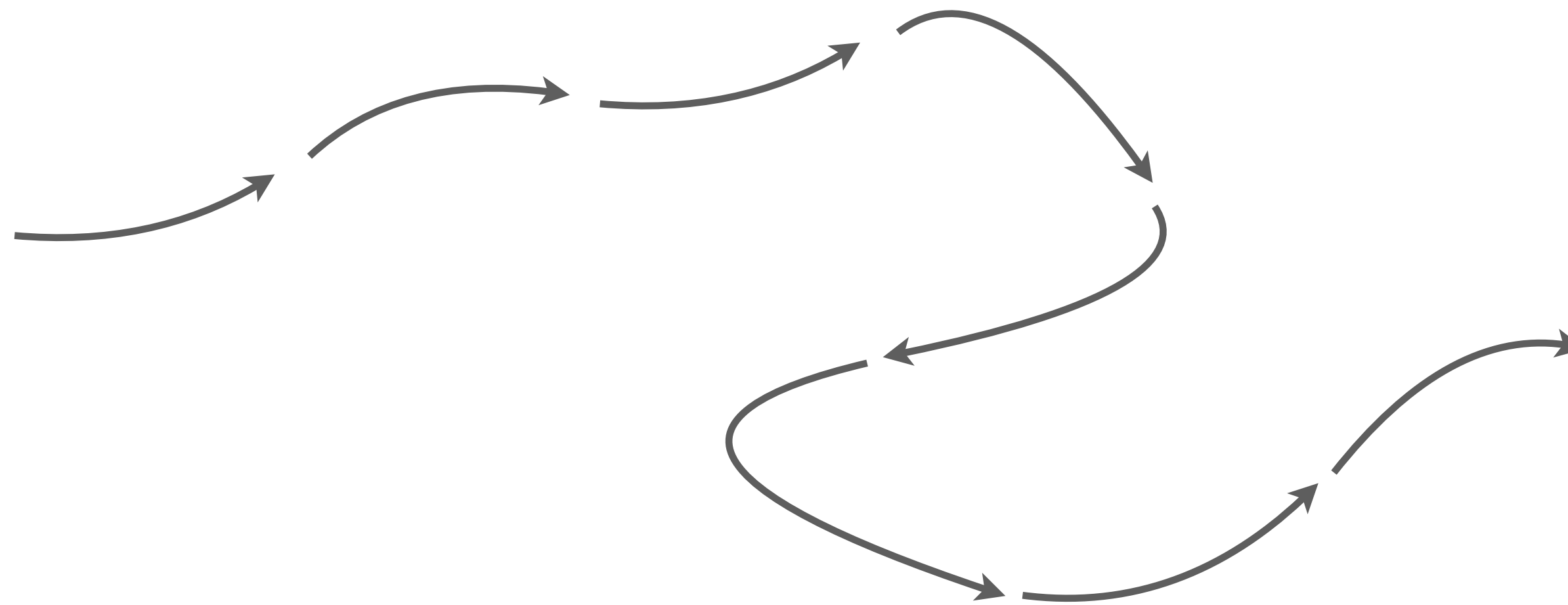
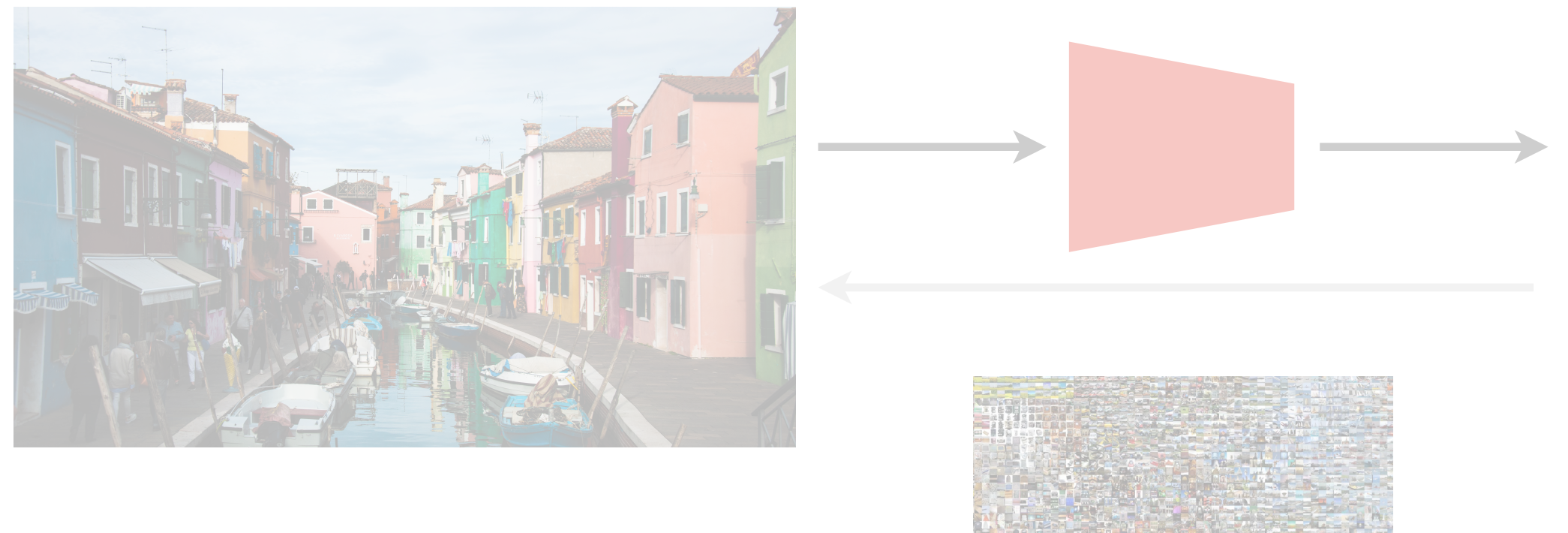
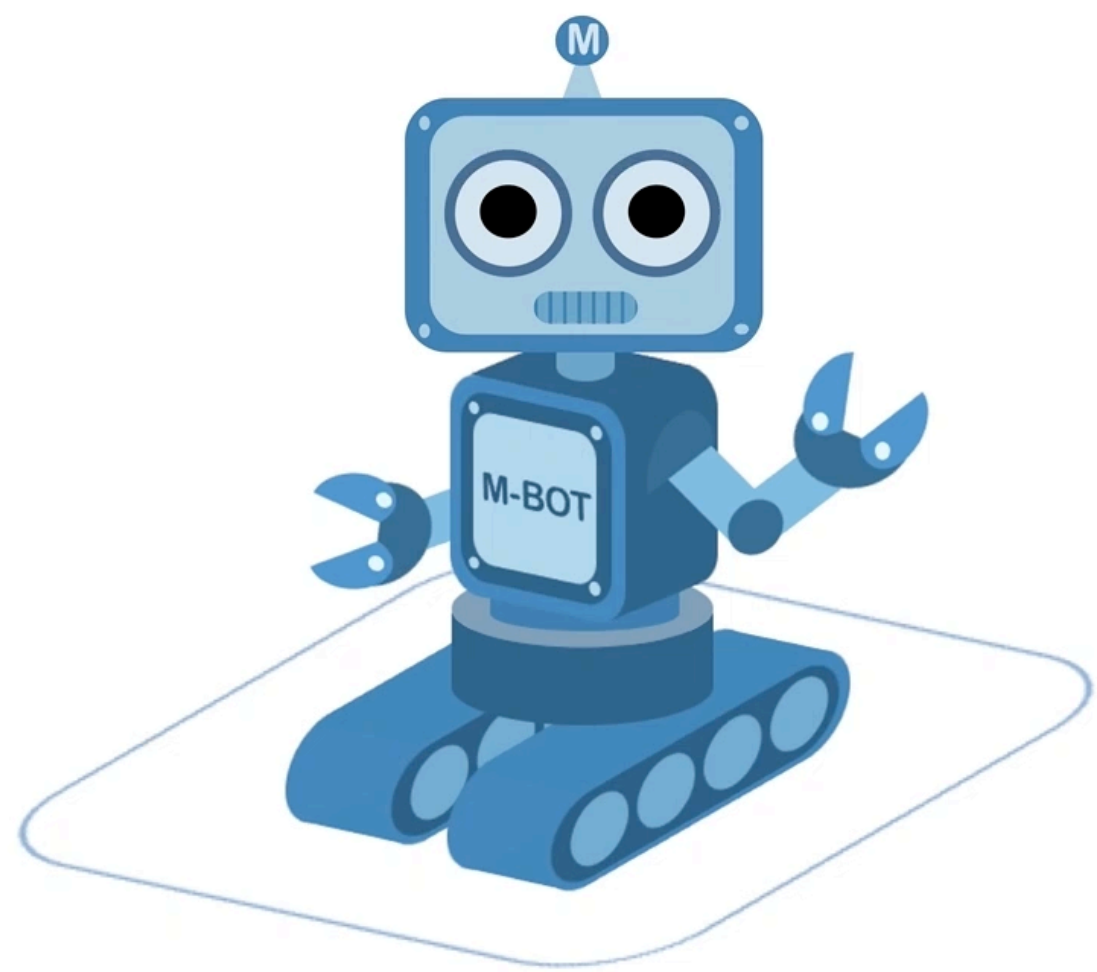
Challenges

Egocentric vision
Active perception
Sparse rewards



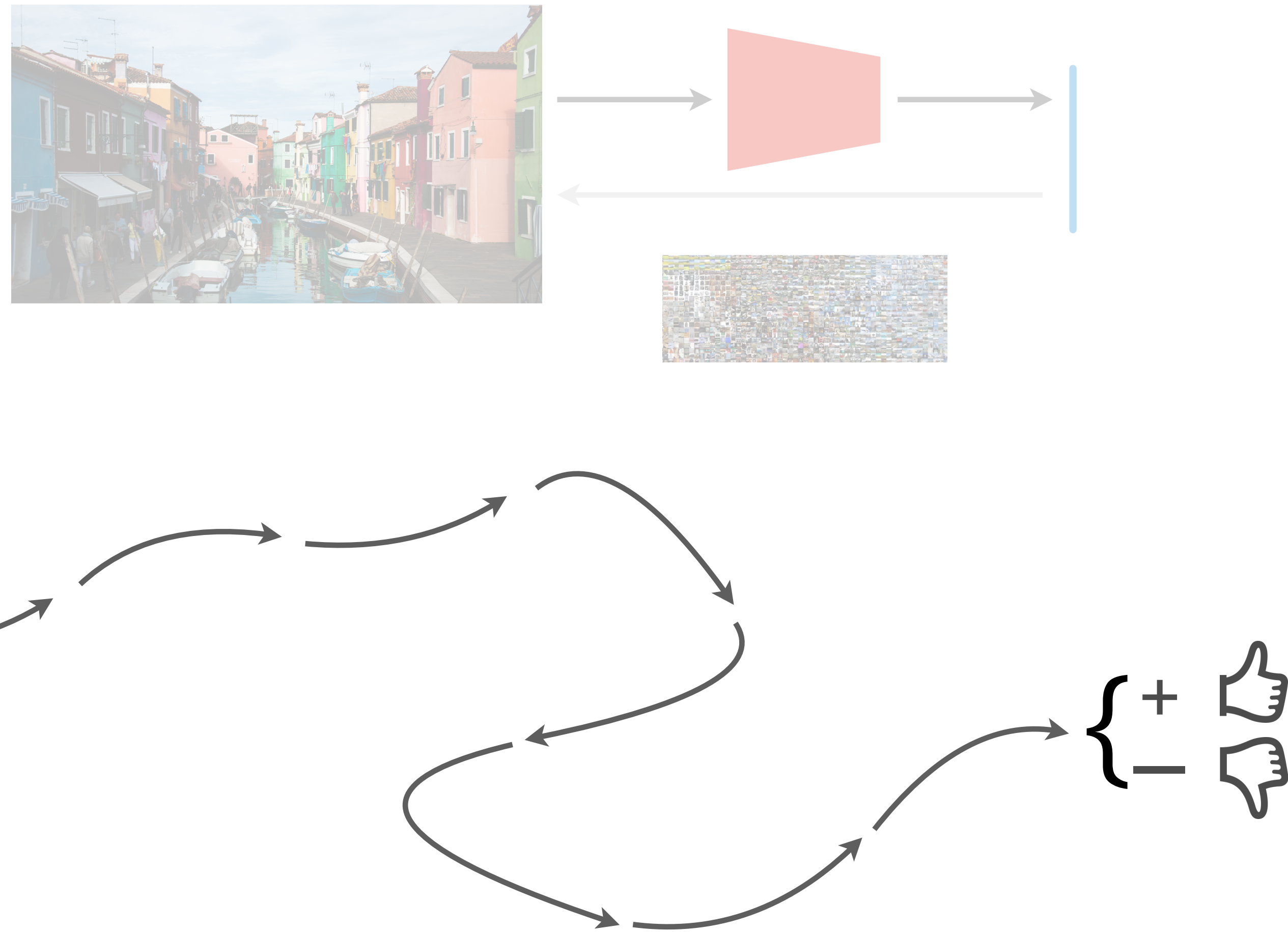
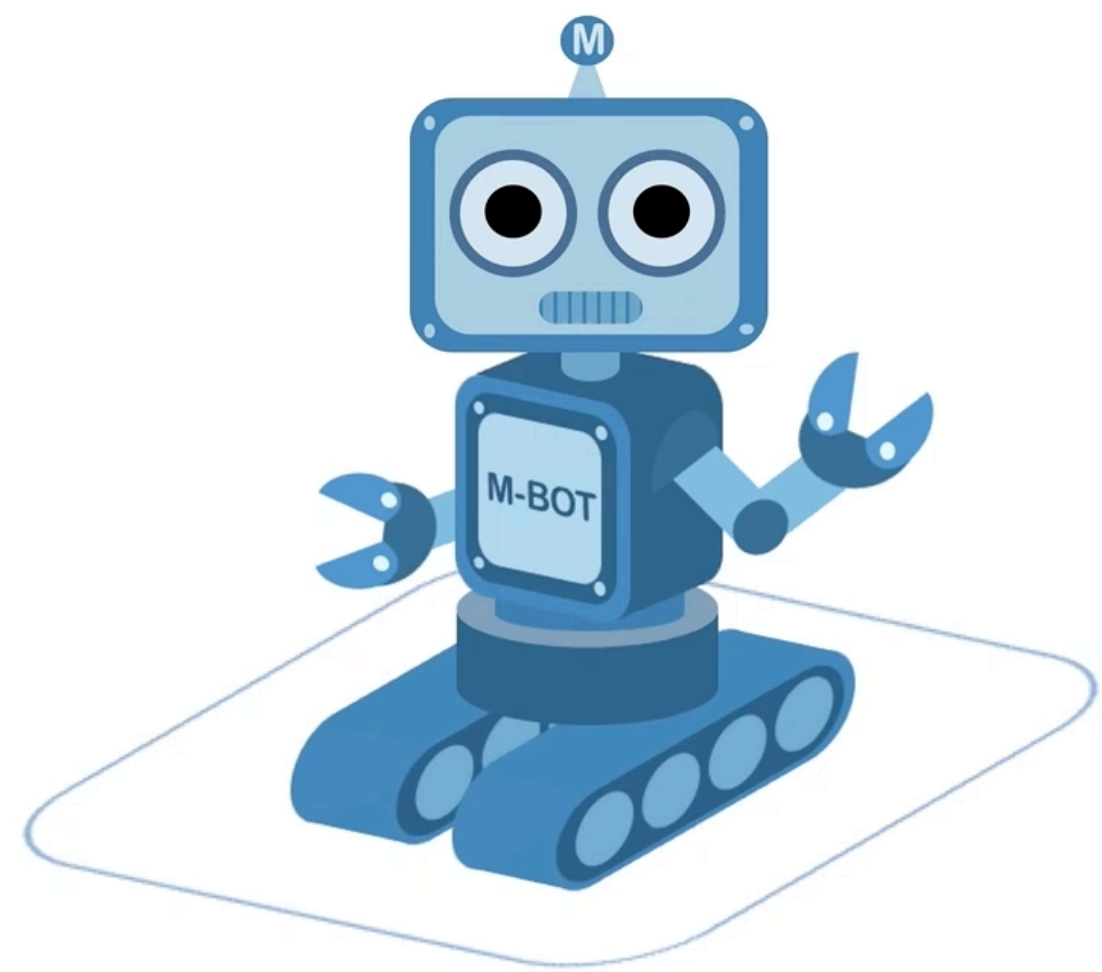
Challenges

Egocentric vision
Active perception
Sparse rewards



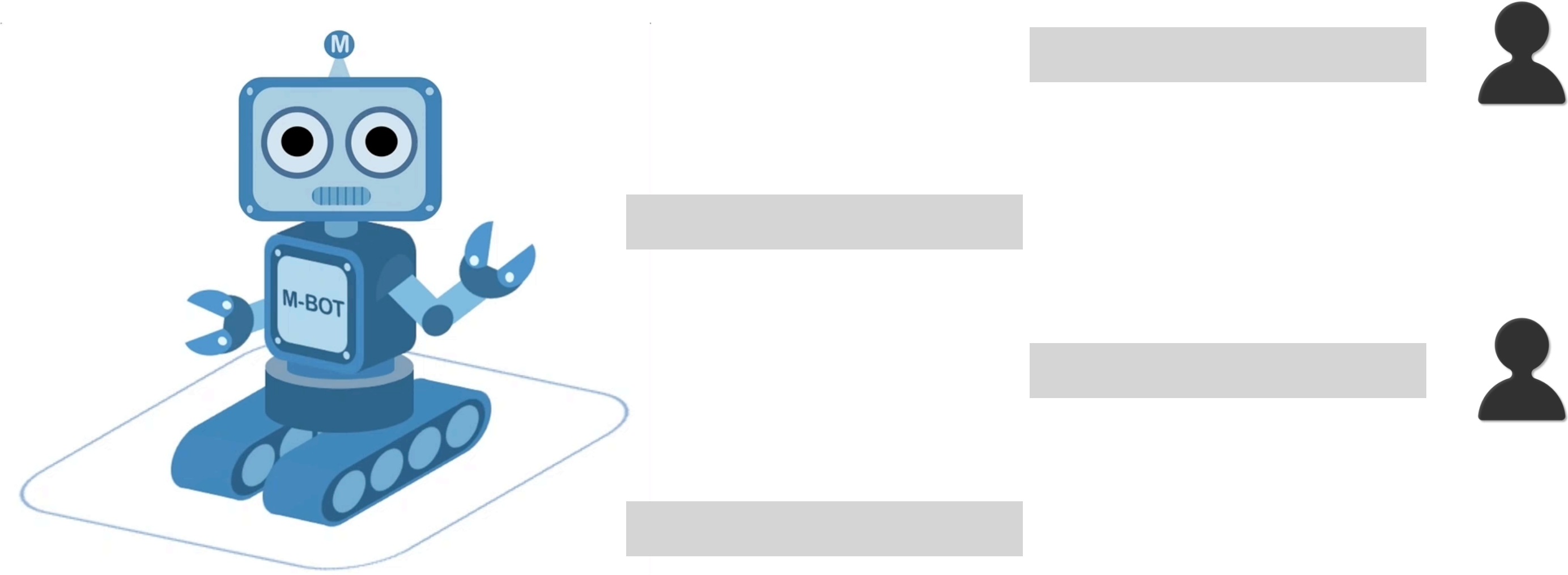
Challenges

Egocentric vision
Active perception
Sparse rewards

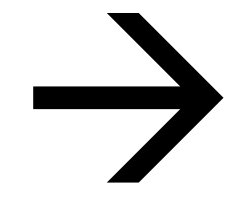


Challenges

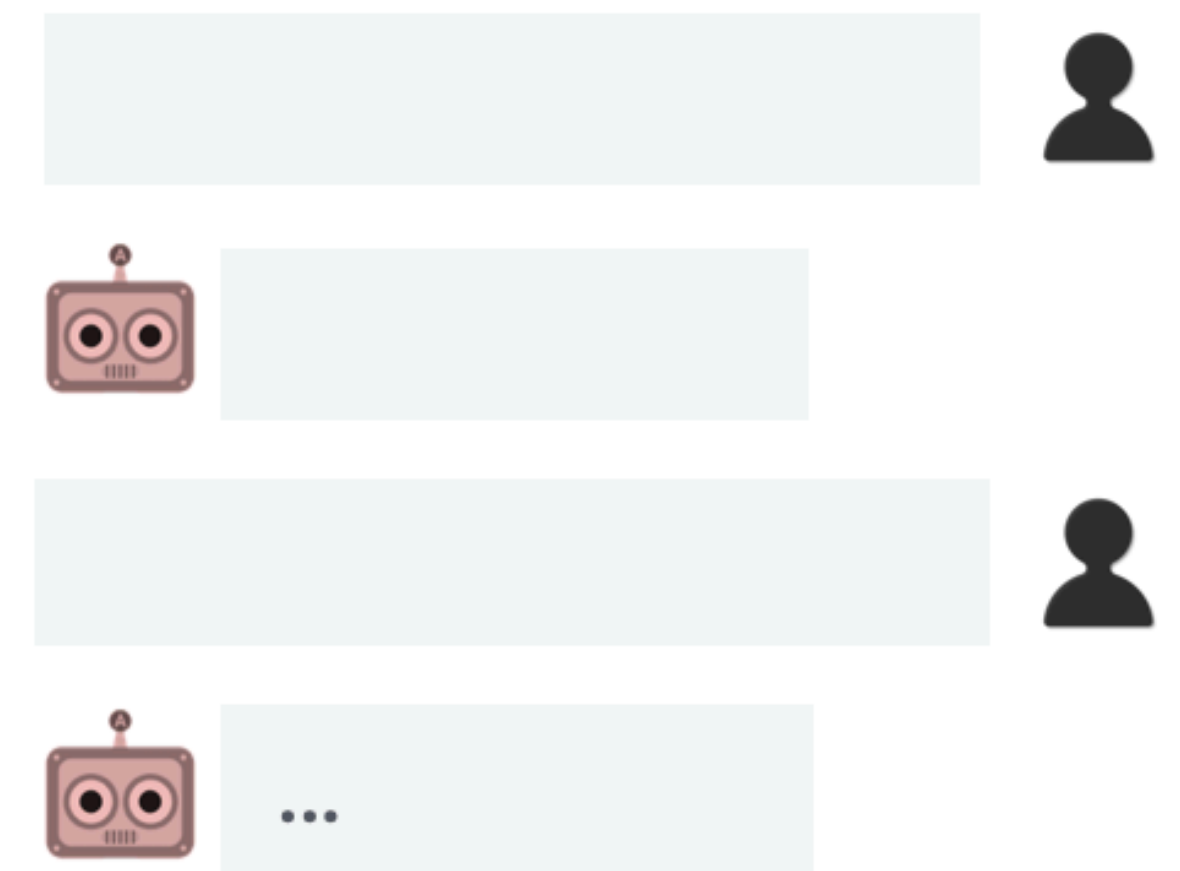
- Egocentric vision
- Active perception
- Sparse rewards
- Language understanding



Internet AI



Embodied AI





Habitat



Manolis Savva^{1,4*}

Abhishek Kadian^{1*}

Oleksandr Maksymets^{1*}

Yili Zhao¹

Erik Wijmans^{1,2,3}

Bhavana Jain¹



Julian Straub²

Jia Liu¹

Vladlen Koltun⁵

Jitendra Malik^{1,6}

Devi Parikh^{1,3}

Dhruv Batra^{1,3}

* denotes equal contribution

facebook
Artificial Intelligence Research

1

facebook
Reality Labs

2

Georgia Tech

3

SFU

4

intel

5

Berkeley
UNIVERSITY OF CALIFORNIA

26

Standardizing the Embodied AI “software stack”

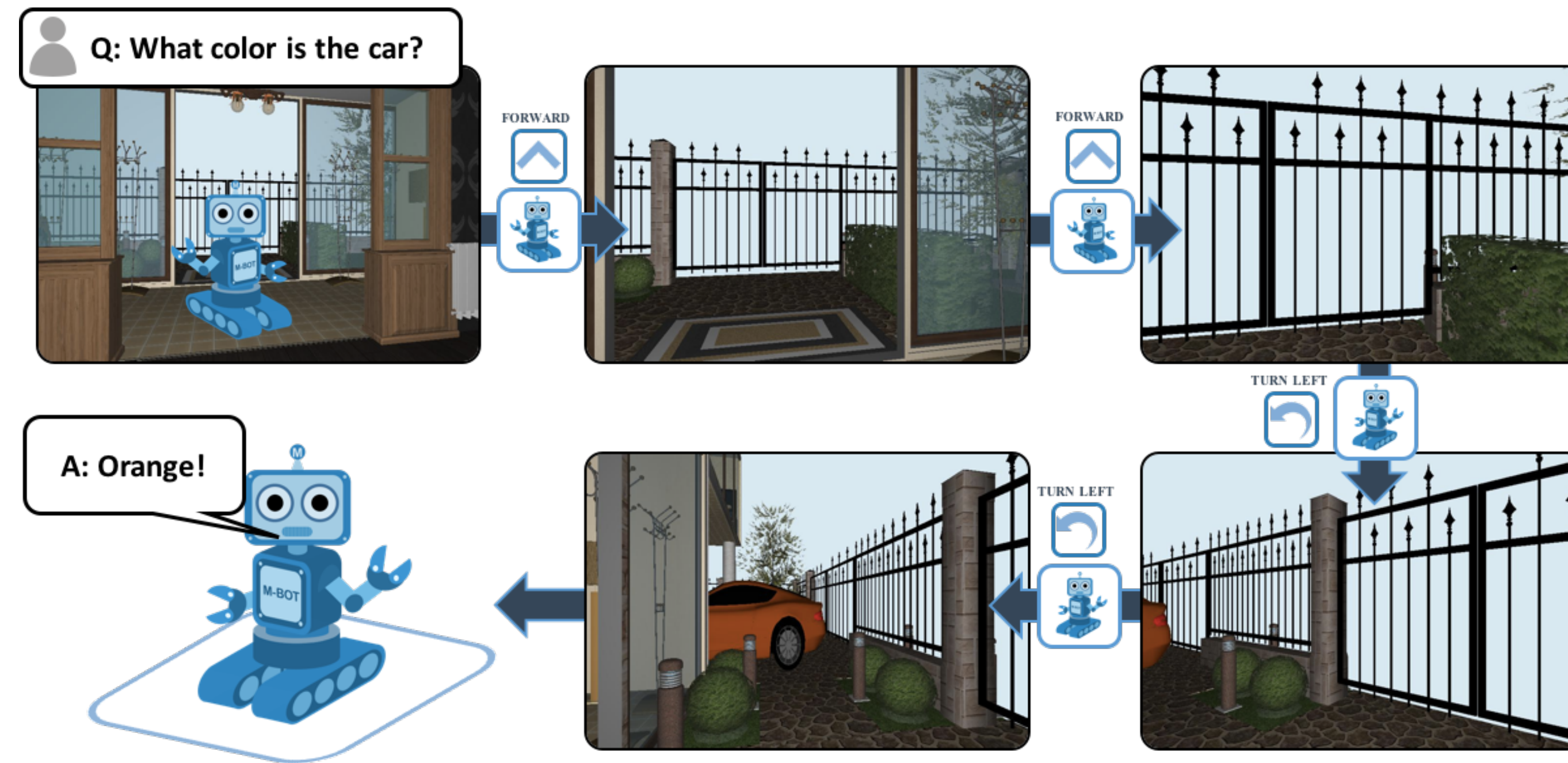
Standardizing the Embodied AI “software stack”

Tasks

The banner displays five task categories for Embodied AI, each with representative images and a title with citation:

- EmbodiedQA** (Das et al., 2018): Shows a robot in a virtual environment with a question "Q: What color is the car?" and an answer "A: Orange!".
- Language grounding** (Hill et al., 2017): Shows a 3D scene with a question "Find object next to the green object" and a score of 113.
- Interactive QA** (Gordon et al., 2018): Shows a grid of images with questions and answers, such as "Q: Is there bread in the room? A: No" and "Q: How many mugs are in the room? A: 3".
- Vision-Language Navigation** (Anderson et al., 2018): Shows a hallway with a goal "Goal: 8.2m" and a blue arrow indicating the path.
- Visual Navigation** (Zhu et al., 2017, Gupta et al., 2017): Shows a 3D scene with three targets labeled "target 1", "target 2", and "target 3".

Standardizing the Embodied AI “software stack”



EmbodiedQA
(Das et al., 2018)

Standardizing the Embodied AI “software stack”

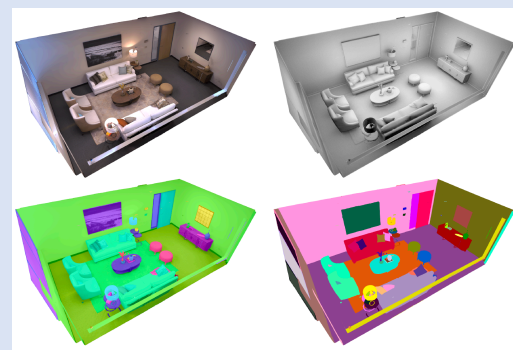


Leave the bedroom, and enter the kitchen. Walk forward, and take a left at the couch. Stop in front of the window.

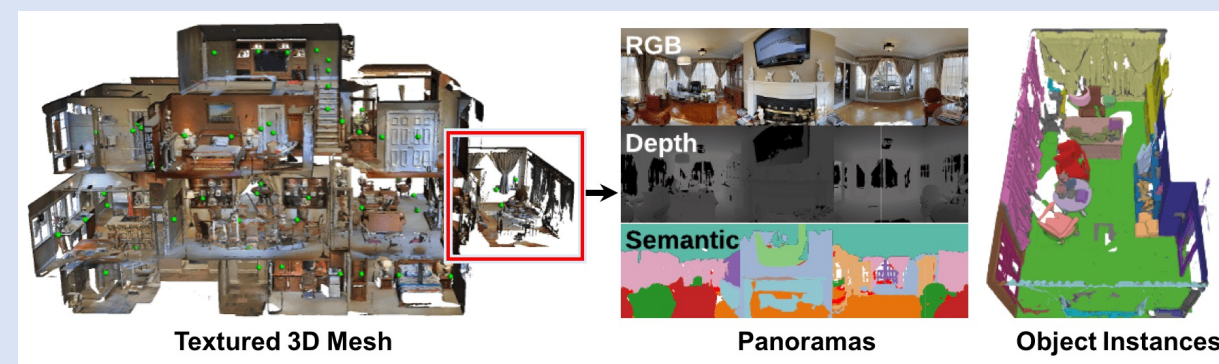
Vision-Language Navigation
(Anderson et al., 2018)

Standardizing the Embodied AI “software stack”

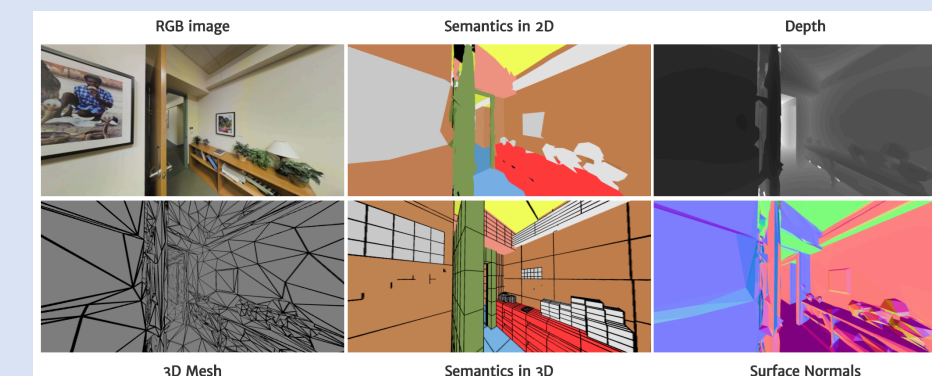
Datasets



Replica (Straub et al., 2019)



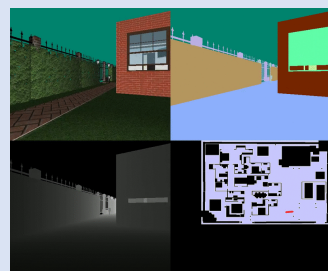
Matterport3D (Chang et al., 2017)



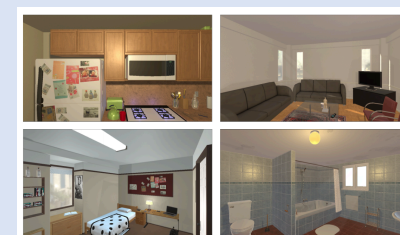
2D-3D-S (Armeni et al., 2017)

Standardizing the Embodied AI “software stack”

Simulators



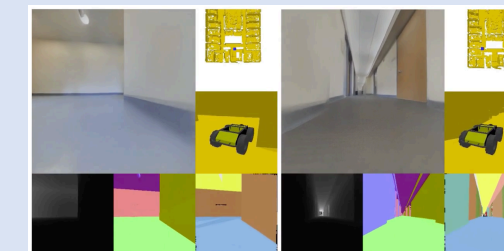
House3D
(Wu et al., 2017)



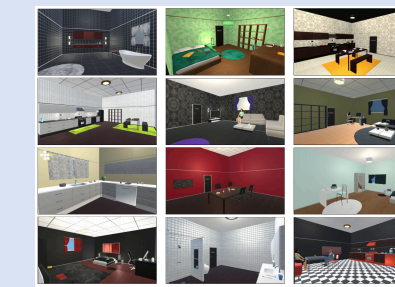
AI2-THOR
(Kolve et al., 2017)



MINOS
(Savva et al., 2017)

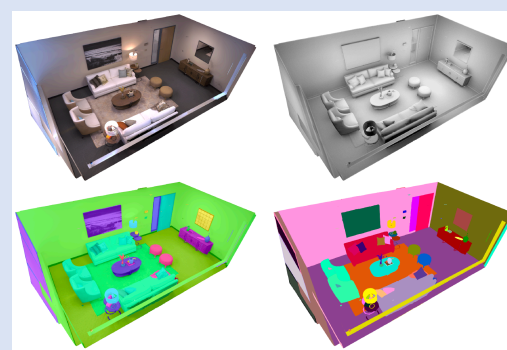


Gibson
(Zamir et al., 2018)

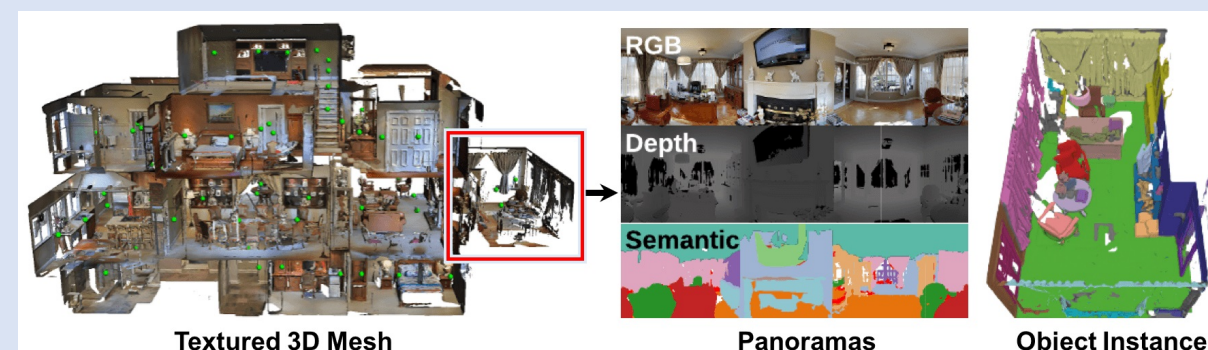


CHALET
(Yan et al., 2018)

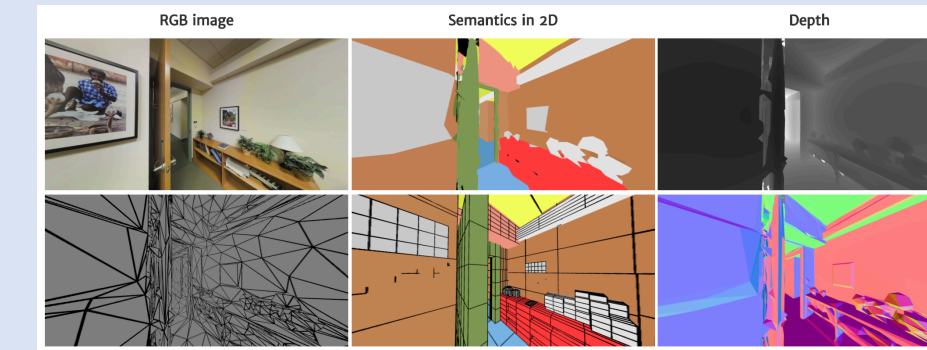
Datasets



Replica (Straub et al., 2019)



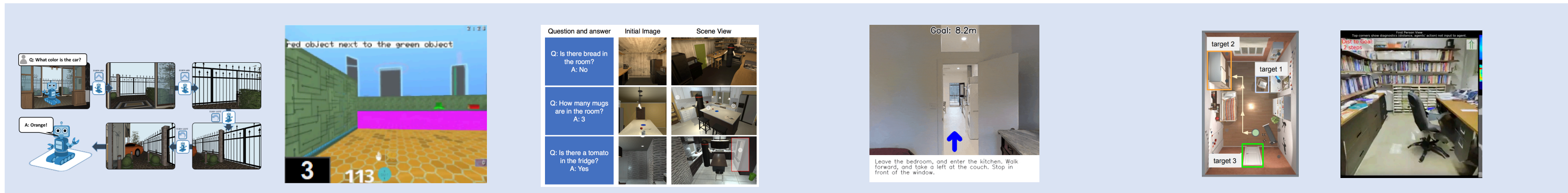
Matterport3D (Chang et al., 2017)



2D-3D-S (Armeni et al., 2017)

Standardizing the Embodied AI “software stack”

Tasks



EmbodiedQA (Das et al., 2018)

Language grounding (Hill et al., 2017)

Interactive QA (Gordon et al., 2018)

Vision-Language Navigation (Anderson et al., 2018)

Visual Navigation (Zhu et al., 2017, Gupta et al., 2017)

Simulators



House3D (Wu et al., 2017)

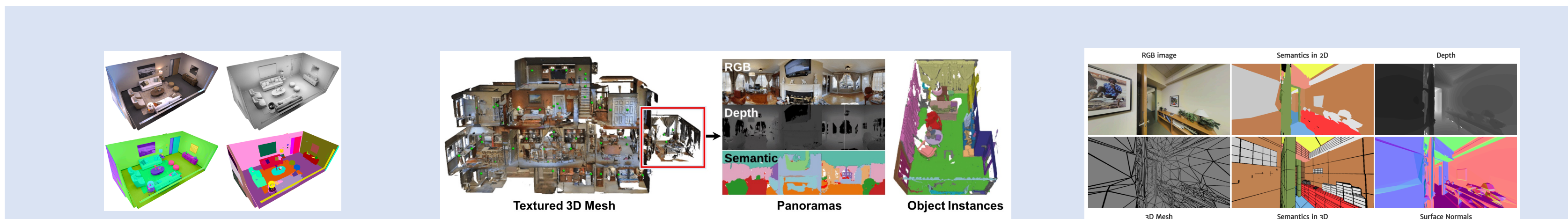
AI2-THOR (Kolve et al., 2017)

MINOS (Savva et al., 2017)

Gibson (Zamir et al., 2018)

CHALET (Yan et al., 2018)

Datasets



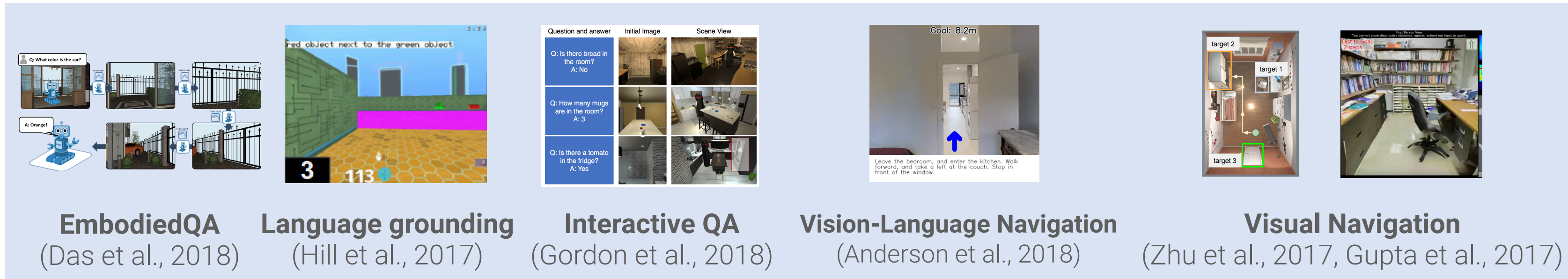
Replica (Straub et al., 2019)

Matterport3D (Chang et al., 2017)

2D-3D-S (Armeni et al., 2017)

Standardizing the Embodied AI “software stack”

Tasks



EmbodiedQA (Das et al., 2018)

Language grounding (Hill et al., 2017)

Interactive QA (Gordon et al., 2018)

Vision-Language Navigation (Anderson et al., 2018)

Visual Navigation (Zhu et al., 2017, Gupta et al., 2017)

Simulators



House3D (Wu et al., 2017)

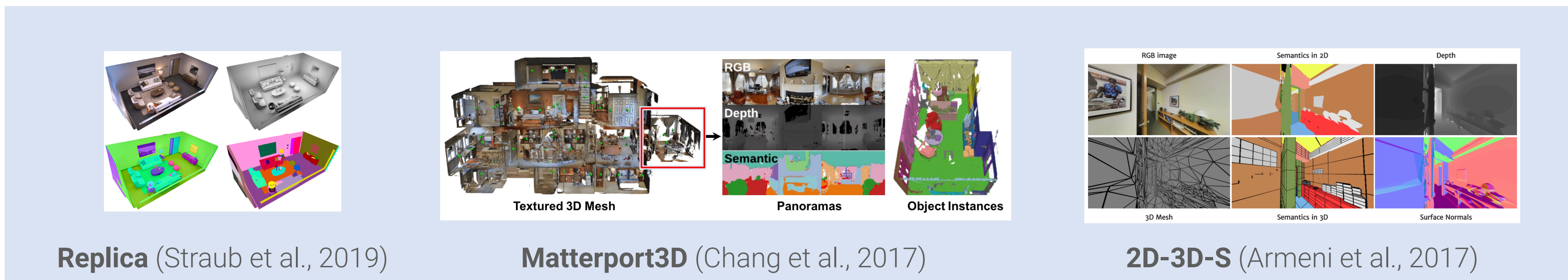
AI2-THOR (Kolve et al., 2017)

MINOS (Savva et al., 2017)

Gibson (Zamir et al., 2018)

CHALET (Yan et al., 2018)

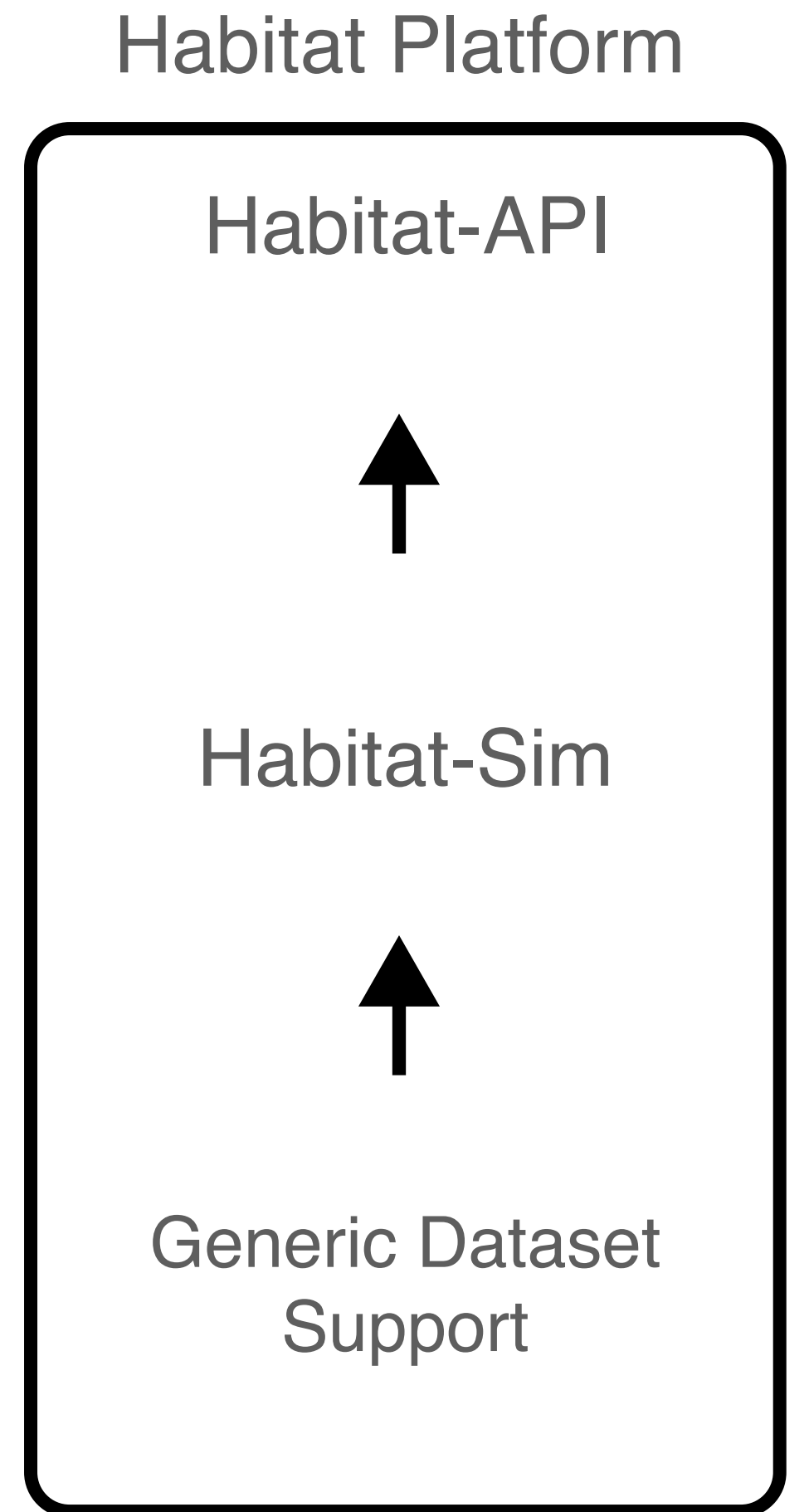
Datasets



Replica (Straub et al., 2019)

Matterport3D (Chang et al., 2017)

2D-3D-S (Armeni et al., 2017)

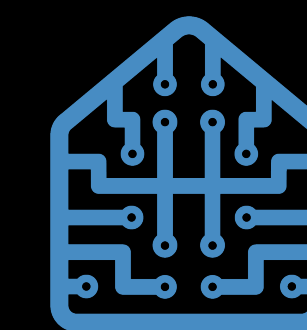


Habitat-Sim Demo

Bring Your Own Scan: Virtualizing Reality



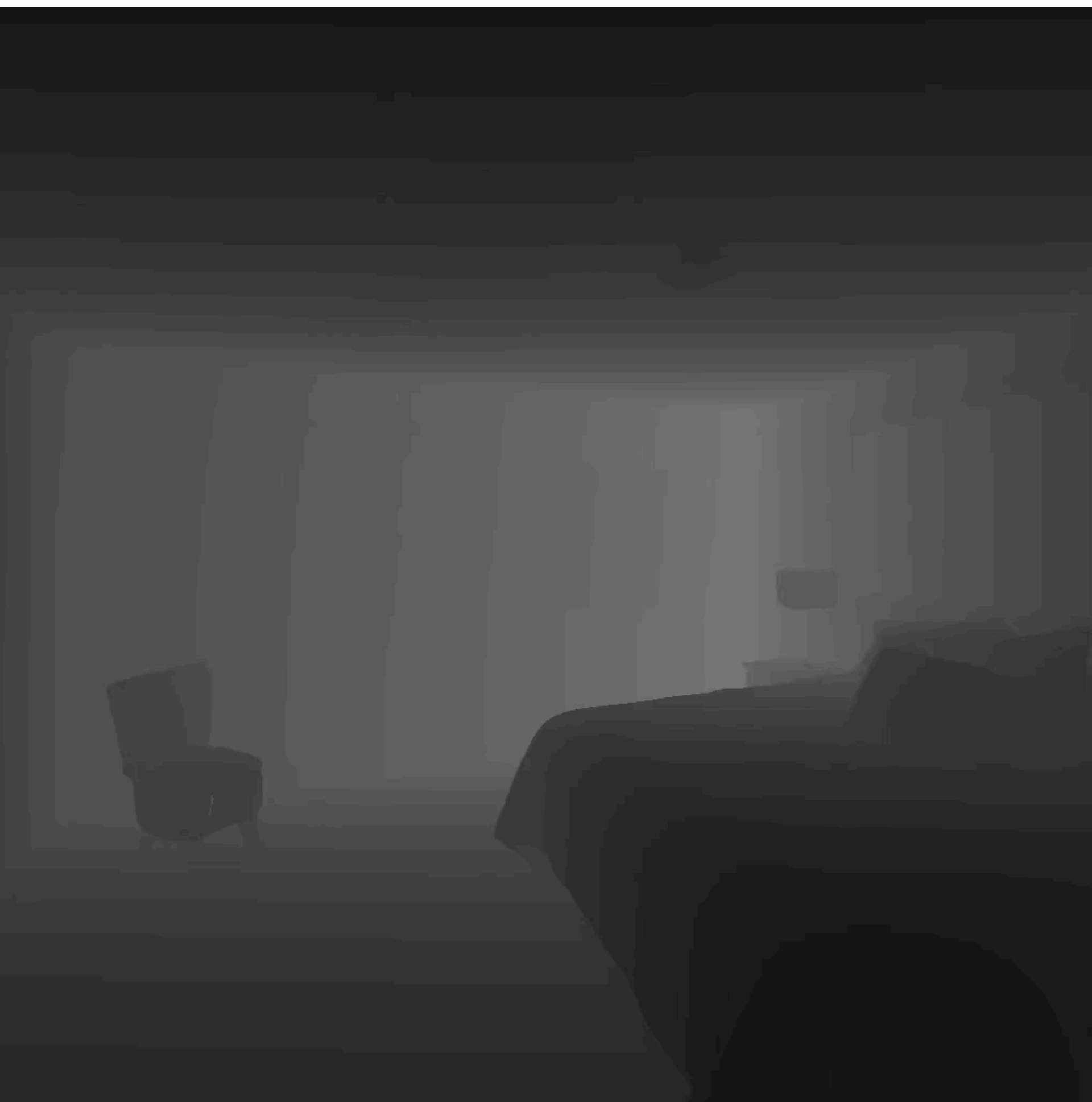
16x



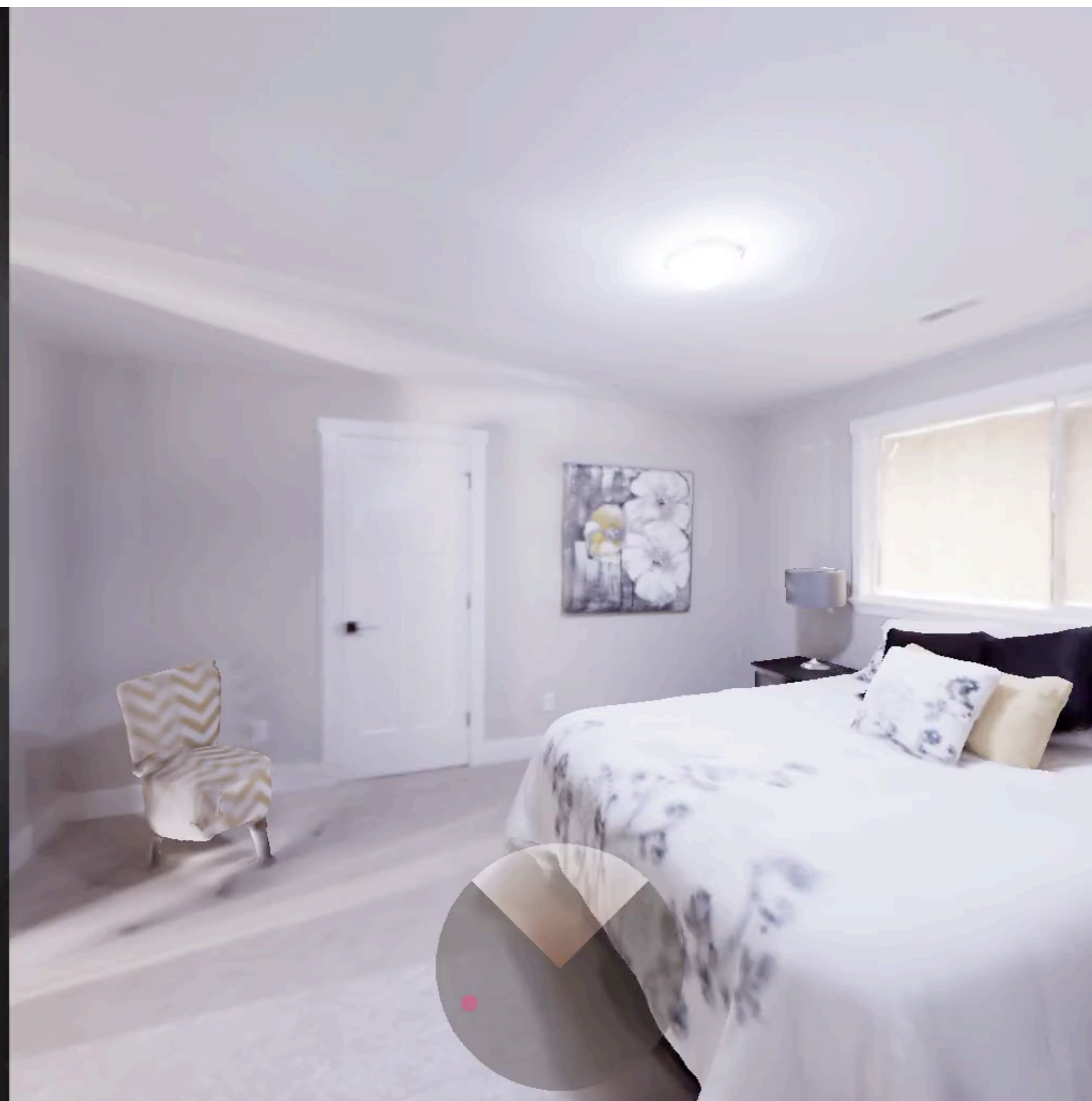
Habitat-API

PointGoal Navigation

PointGoal Navigation



Depth

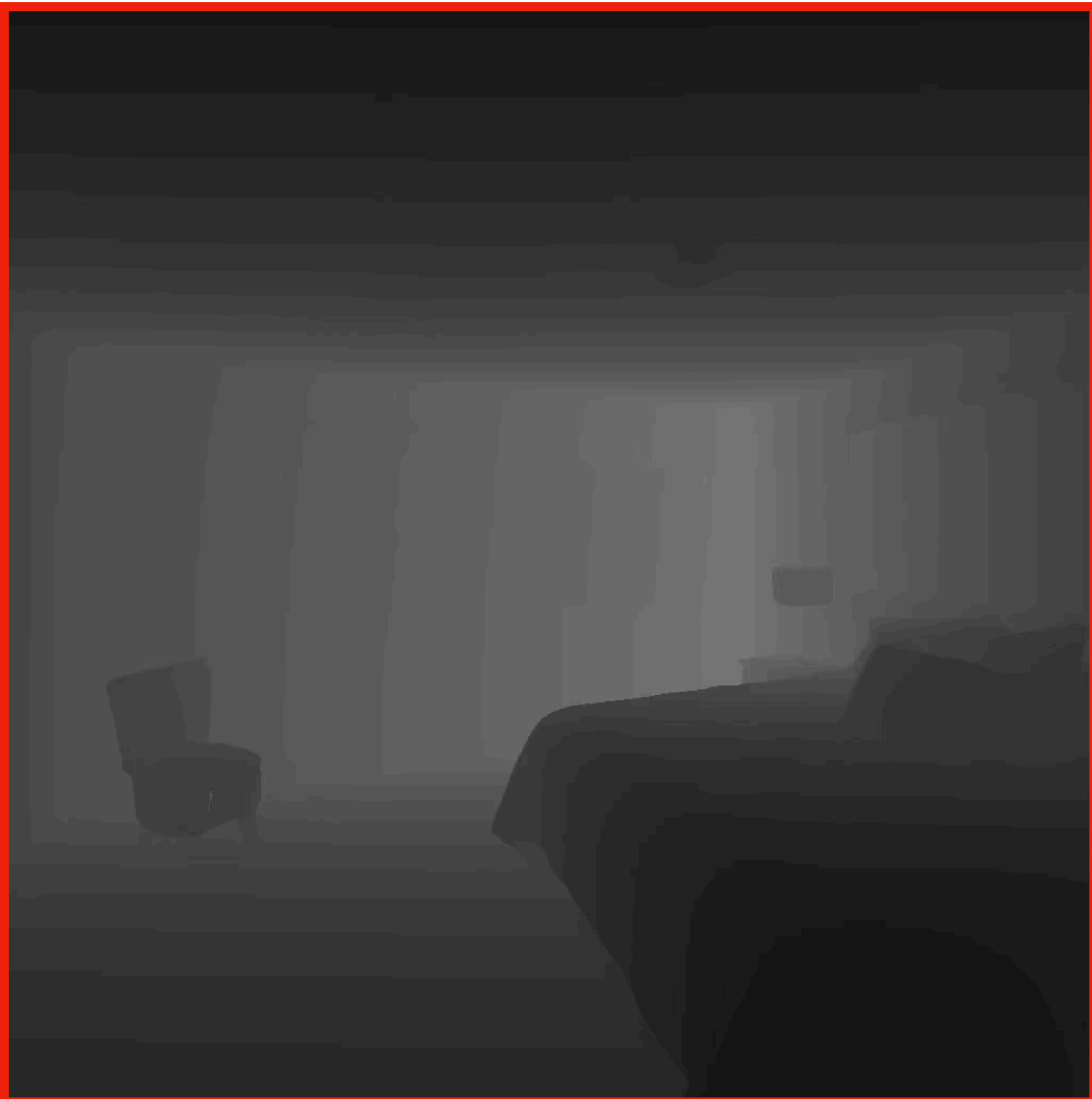


RGB and GPS+Compass



Top Down Map

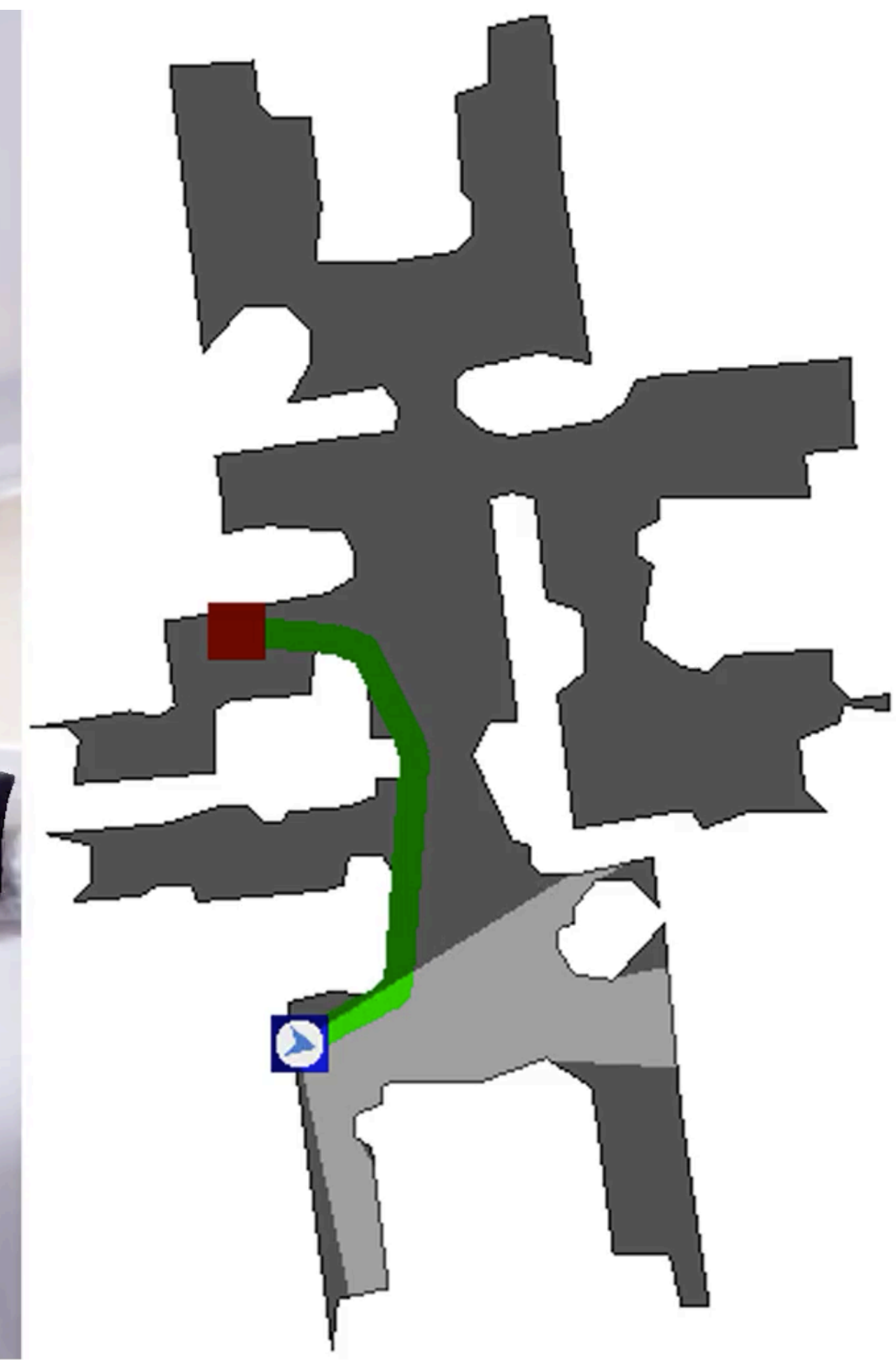
PointGoal Navigation



Depth

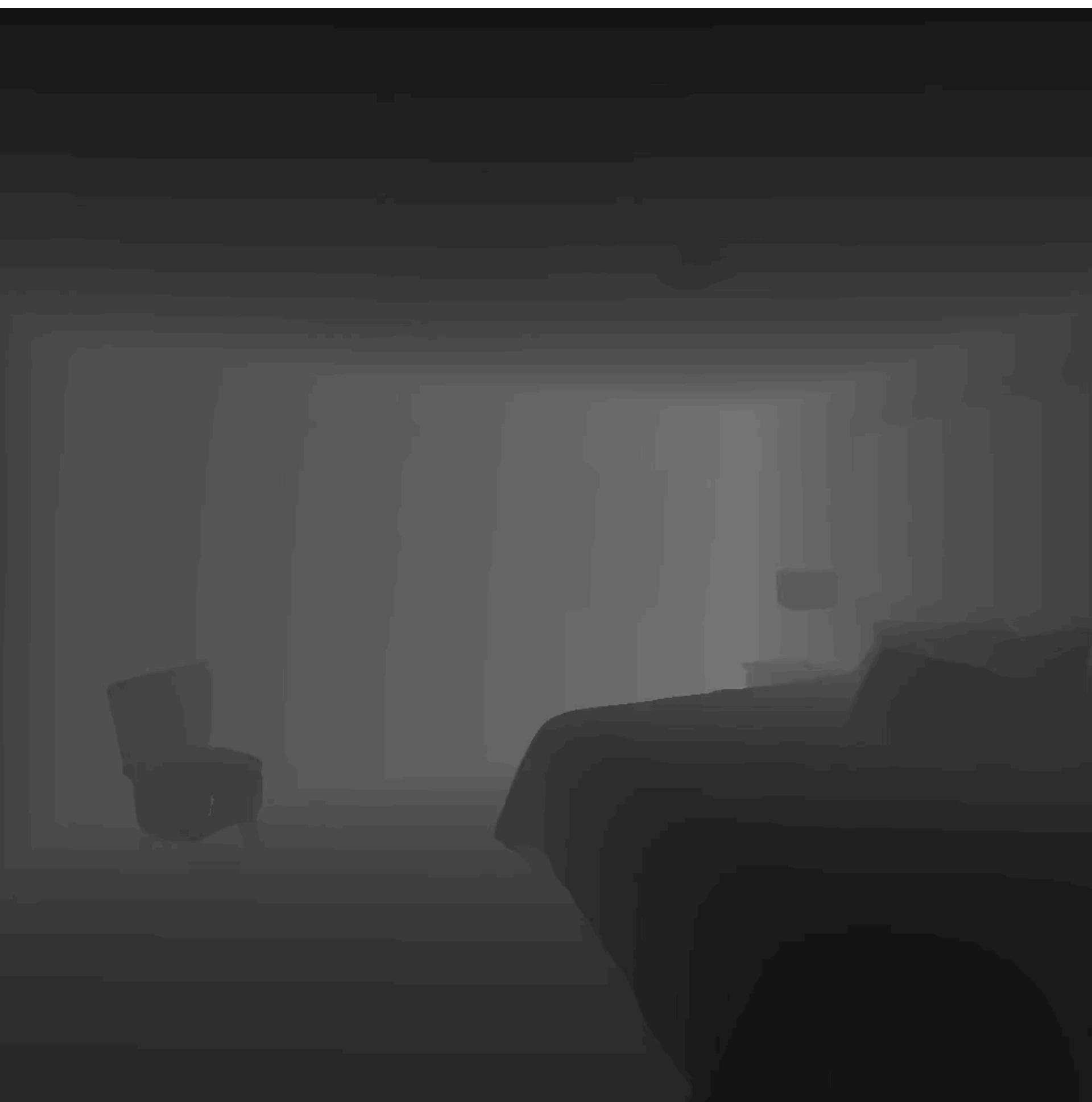


RGB and GPS+Compass

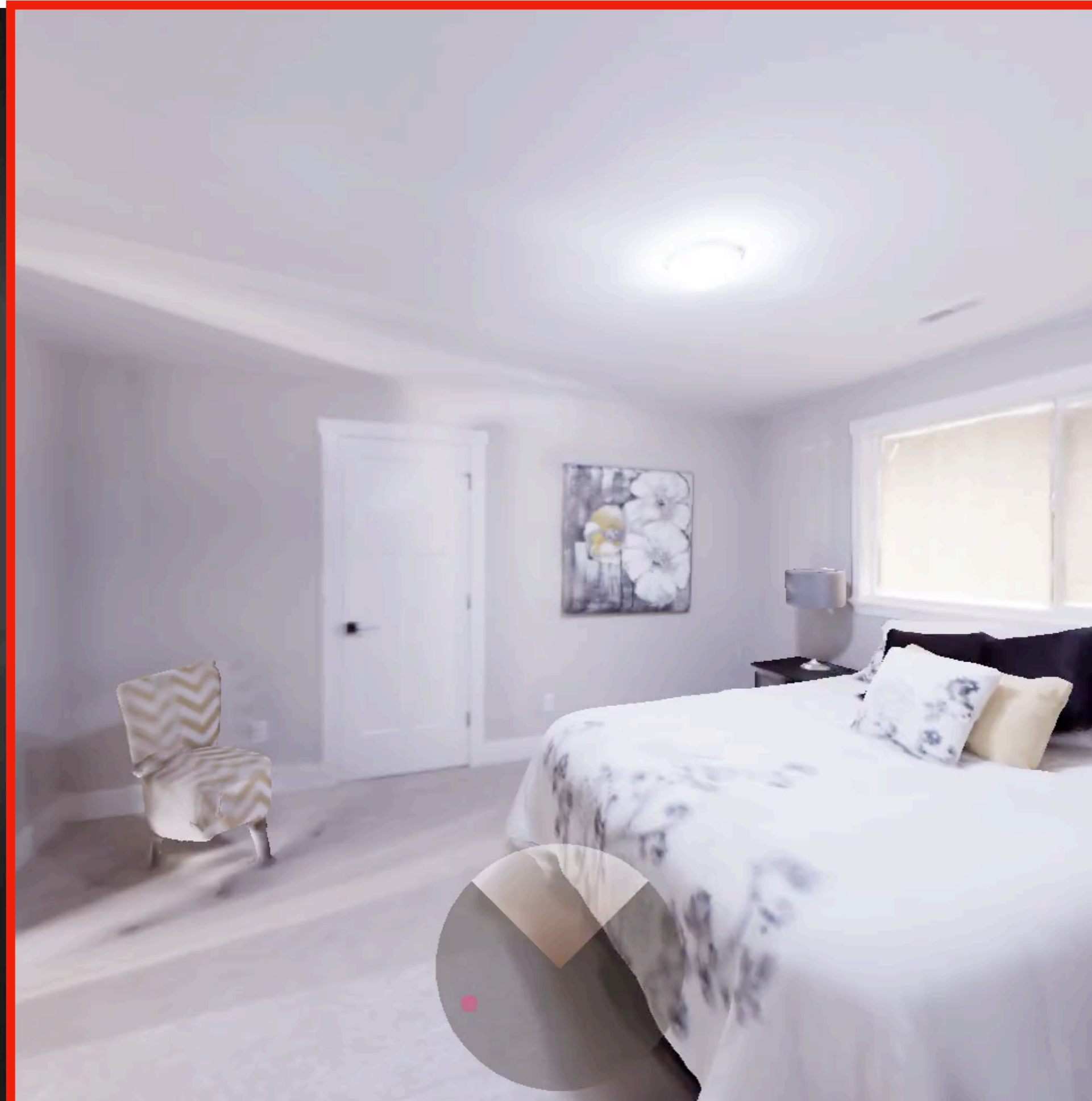


Top Down Map

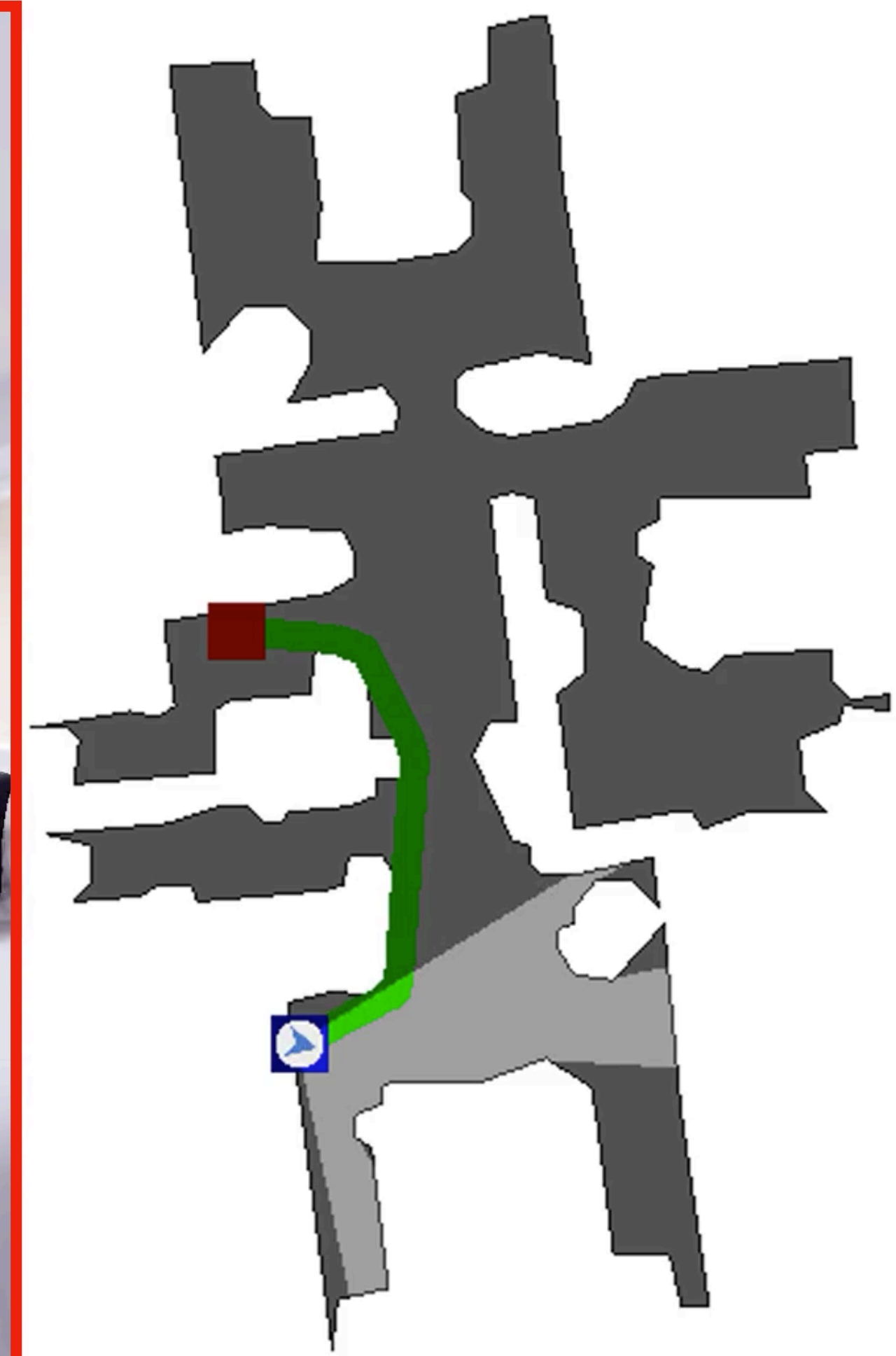
PointGoal Navigation



Depth

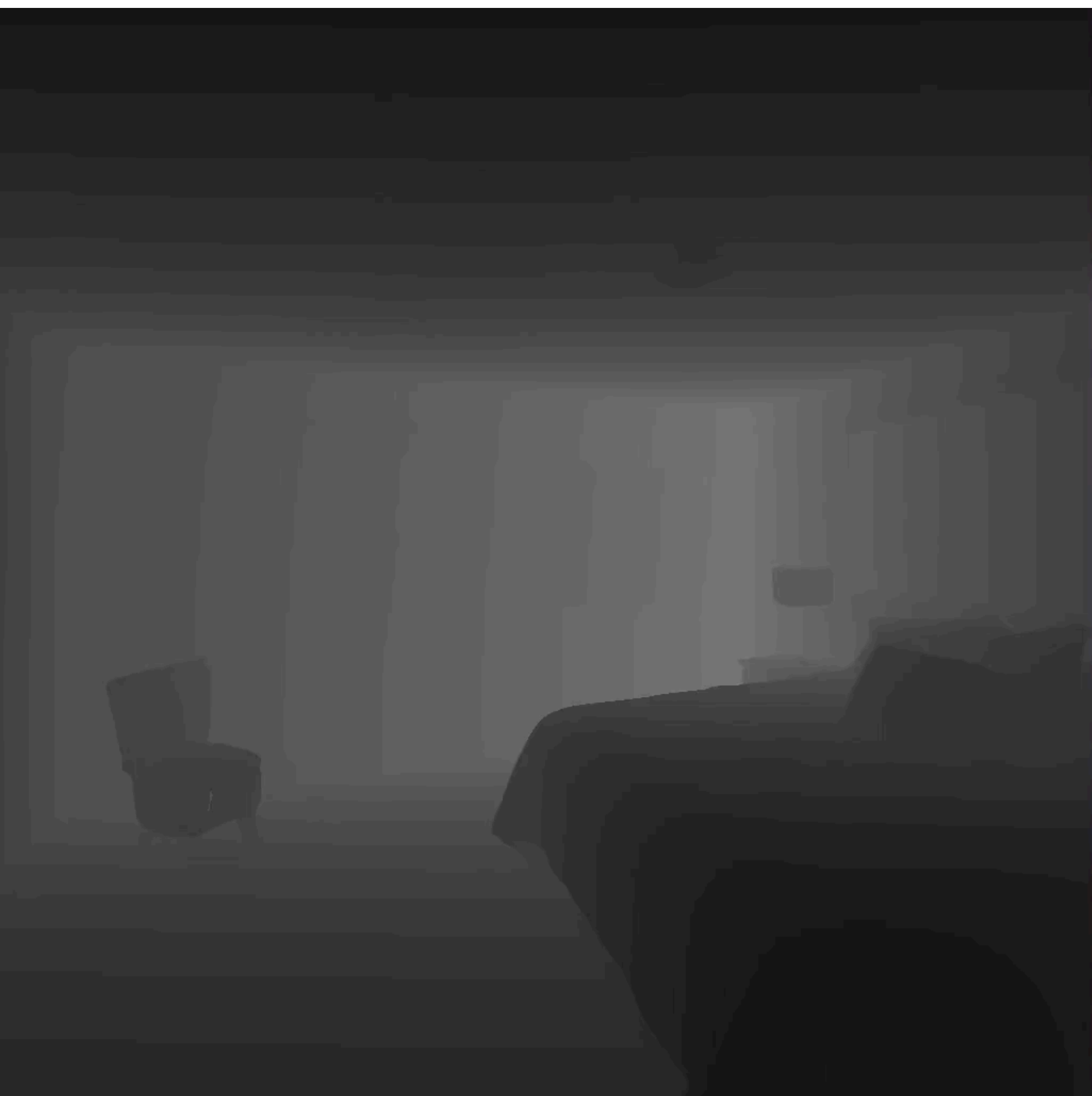


RGB and GPS+Compass

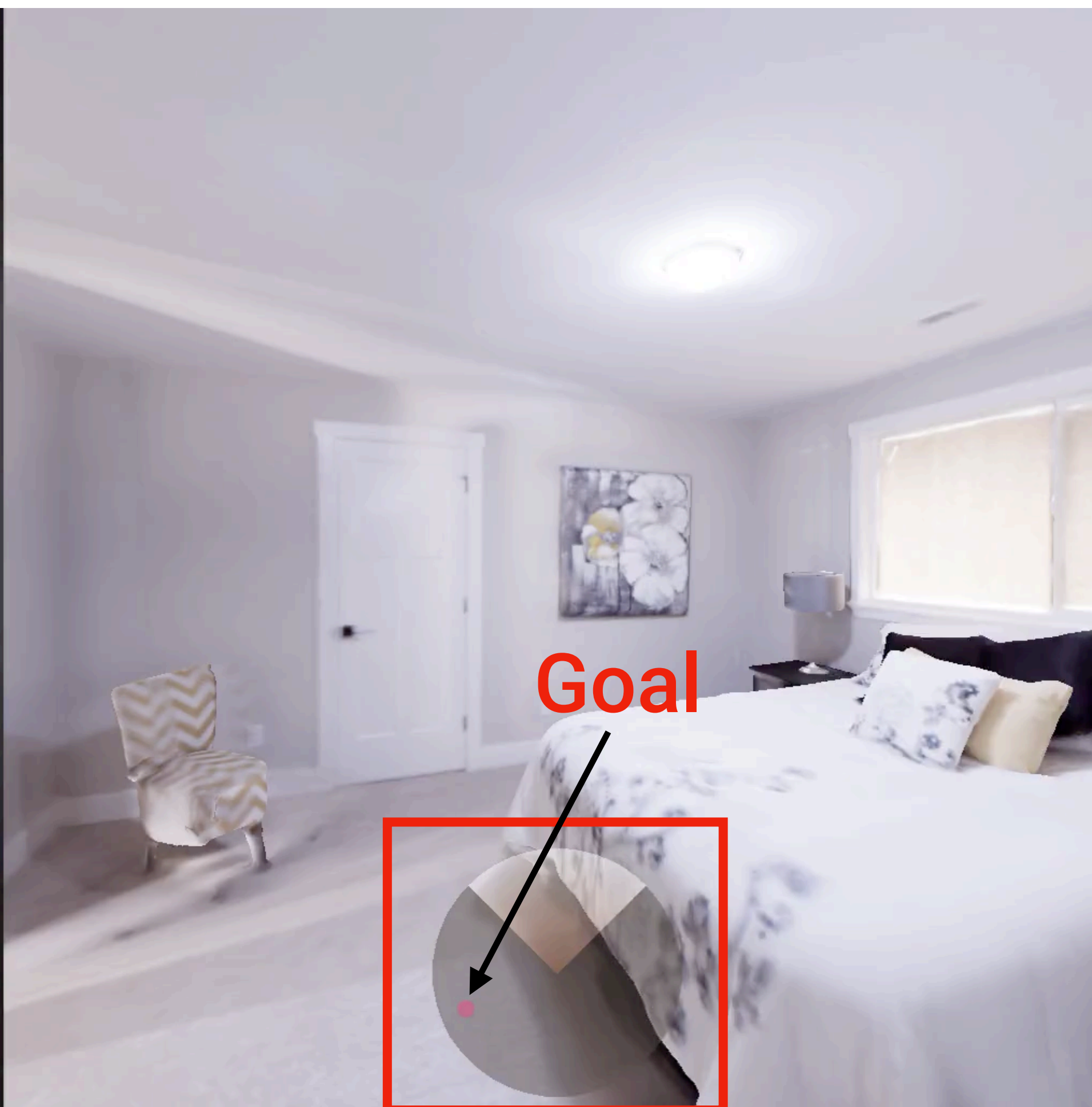


Top Down Map

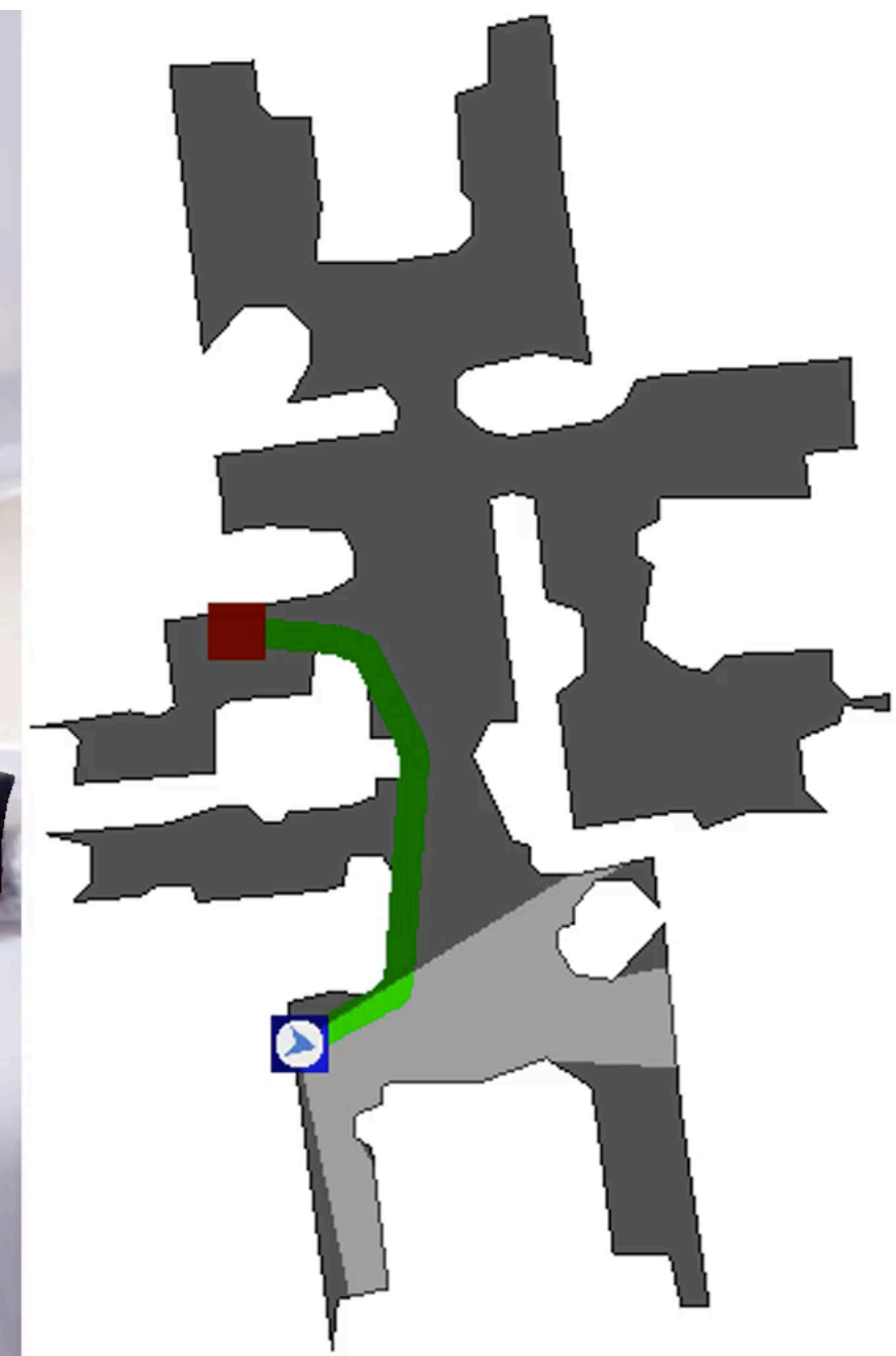
PointGoal Navigation



Depth

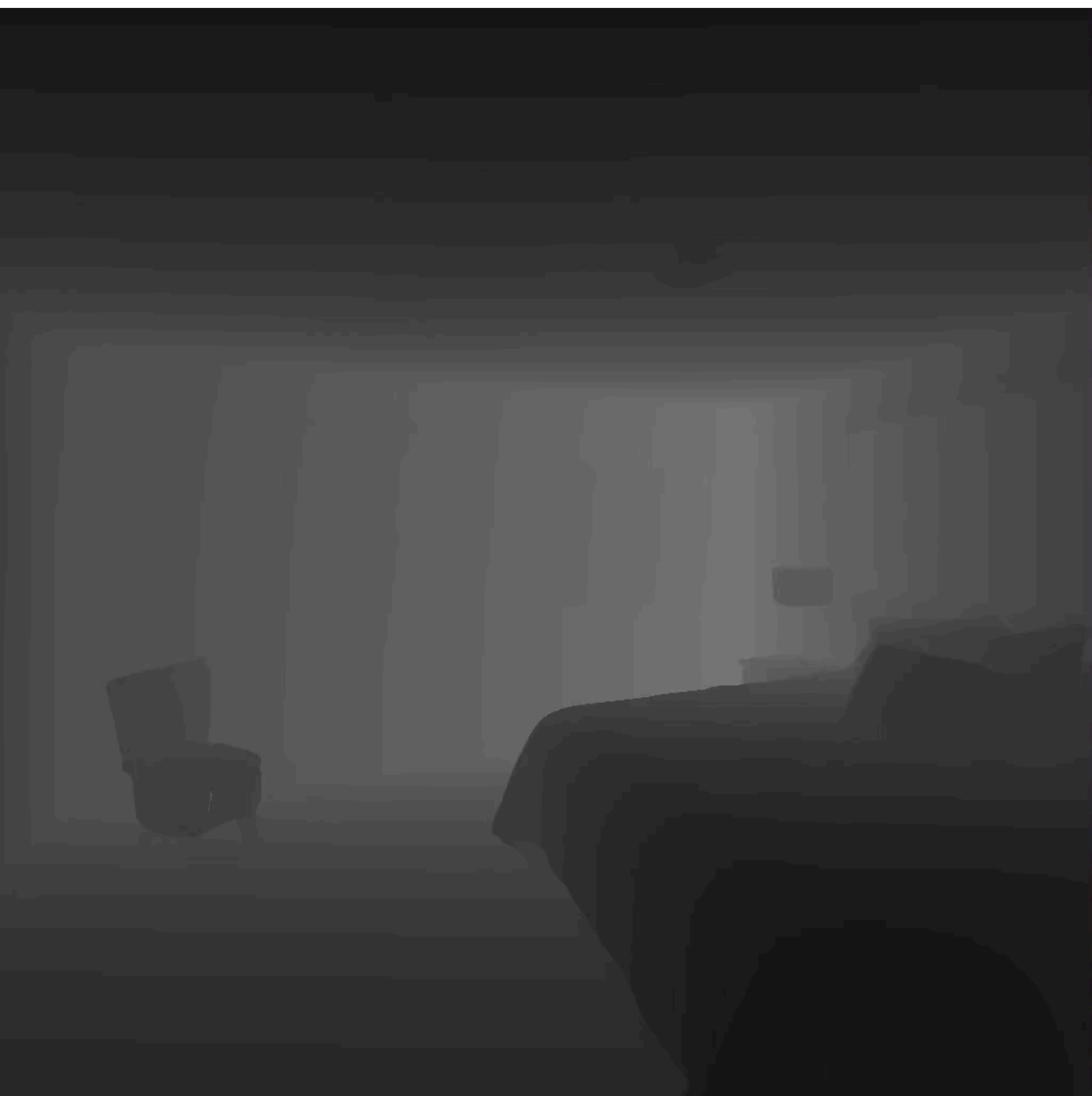


RGB and GPS+Compass

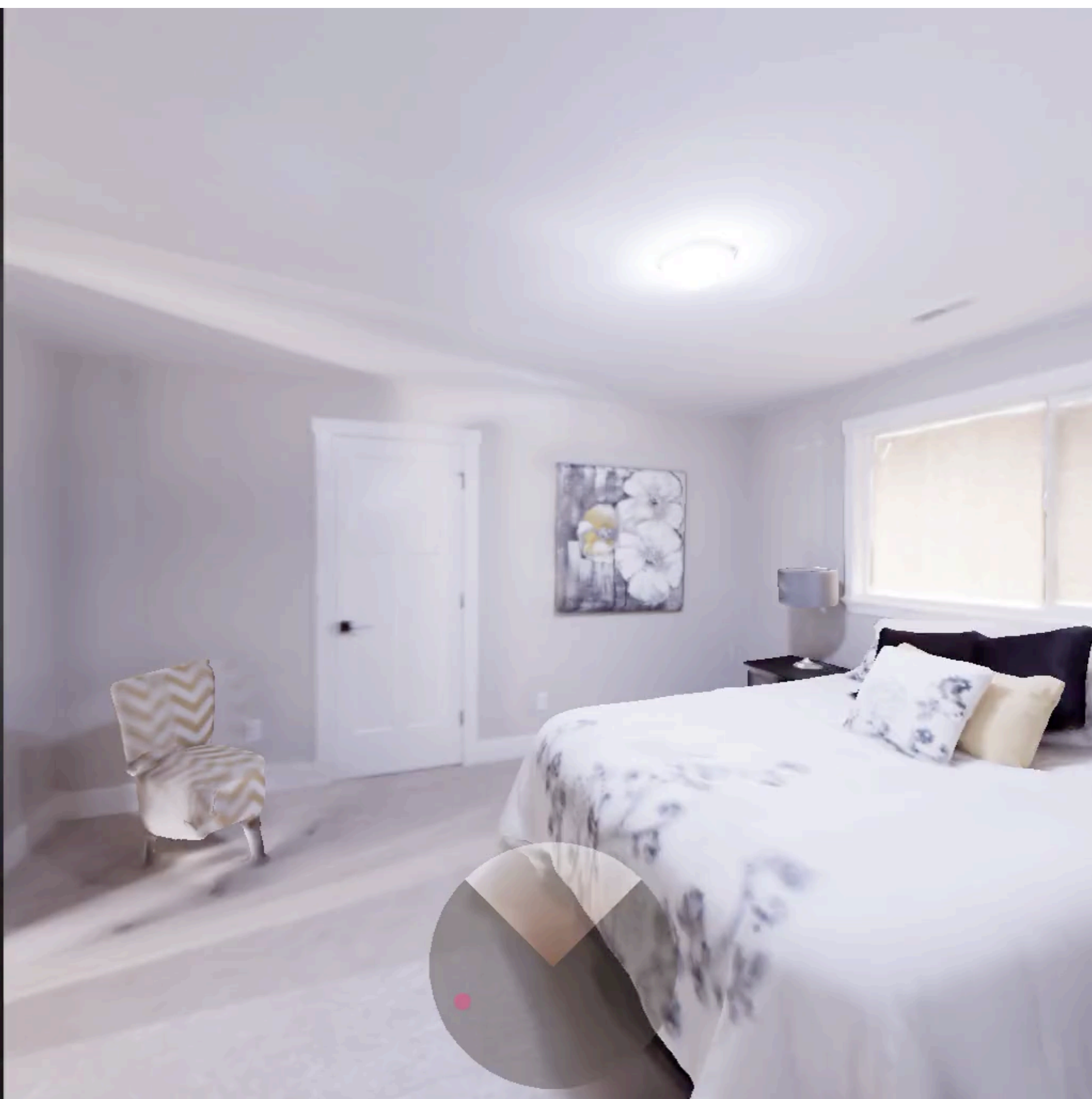


Top Down Map

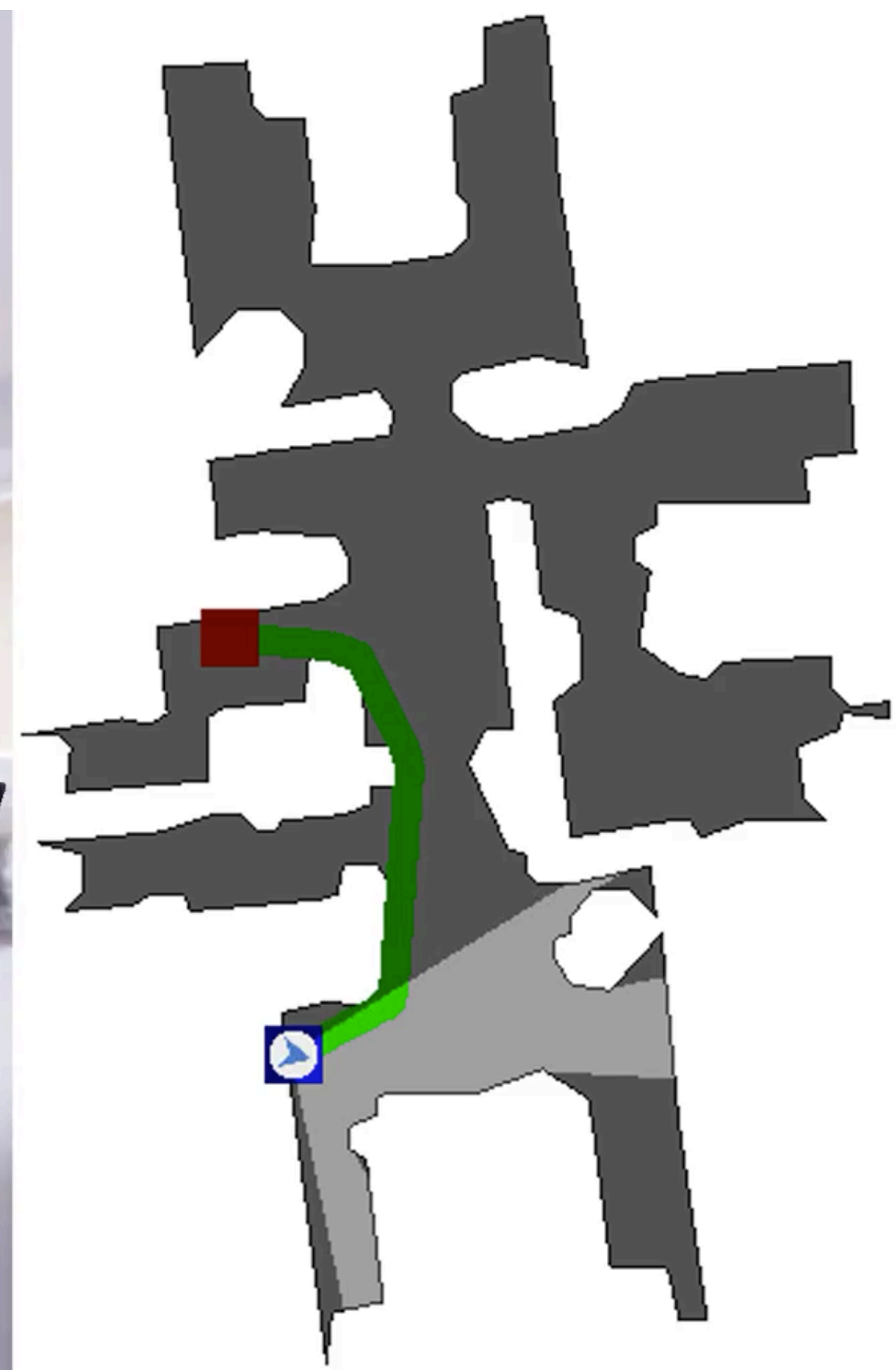
PointGoal Navigation



Depth



RGB and GPS+Compass



Top Down Map

Agent and Model Design

Agent and Model Design



Agent and Model Design



Agent and Model Design



- 0.6m tall cylinder with 0.17m radius

Agent and Model Design



- 0.6m tall cylinder with 0.2m radius
- Actions:
 - <stop>: Indicates the agent believes it has completed the task
 - <forward>: Moves 0.25m forward
 - <left>, <right>: Turn 30 degrees

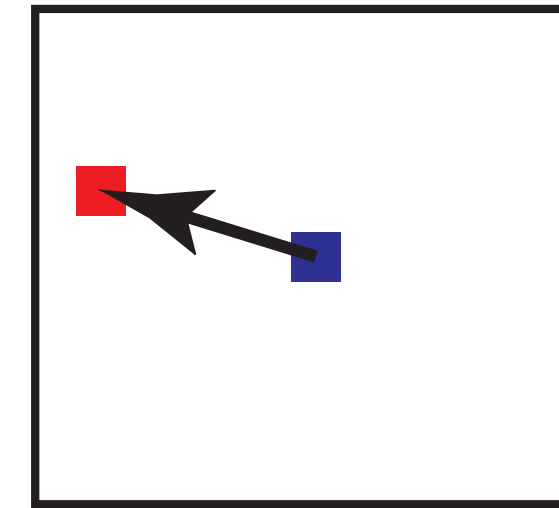
Agent and Model Design



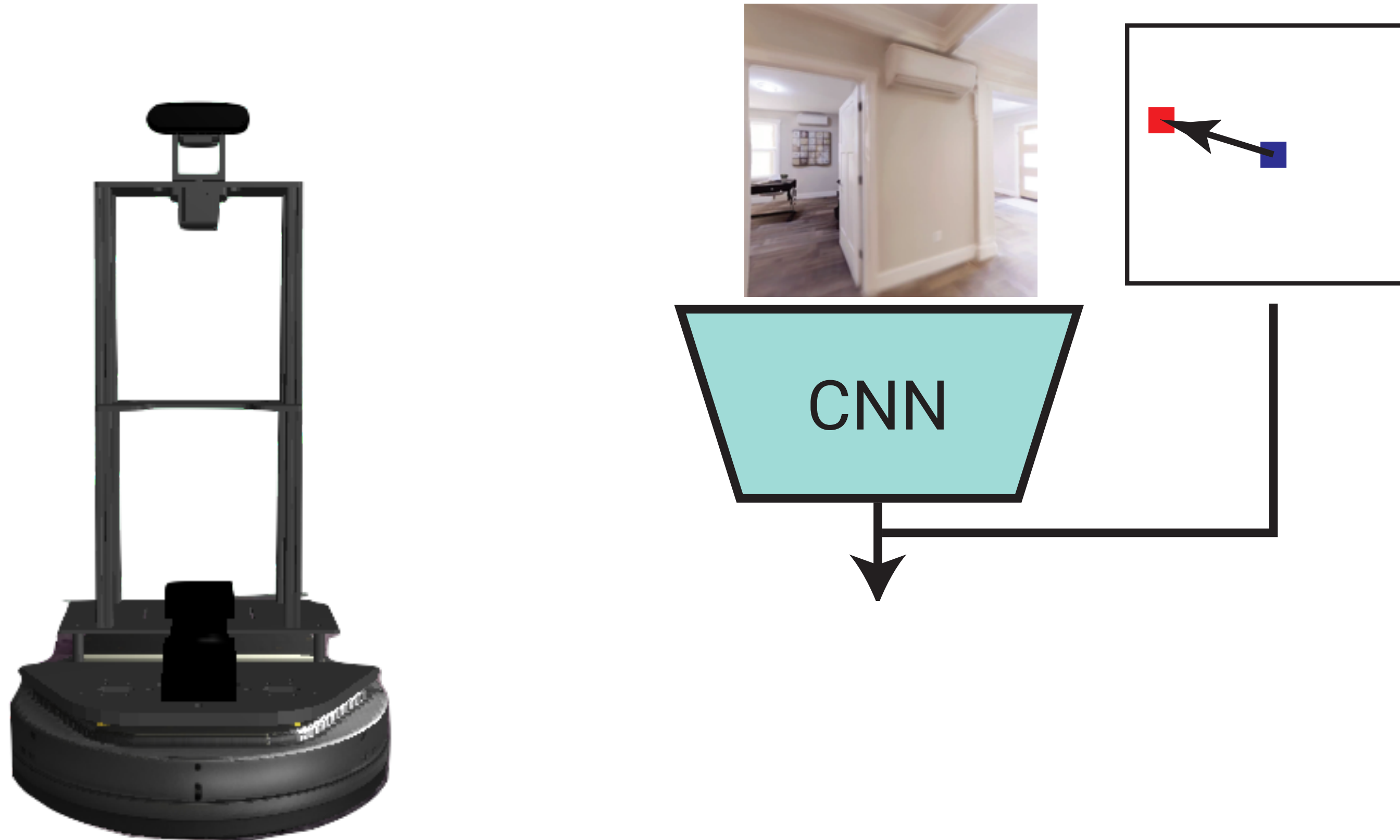
Agent and Model Design



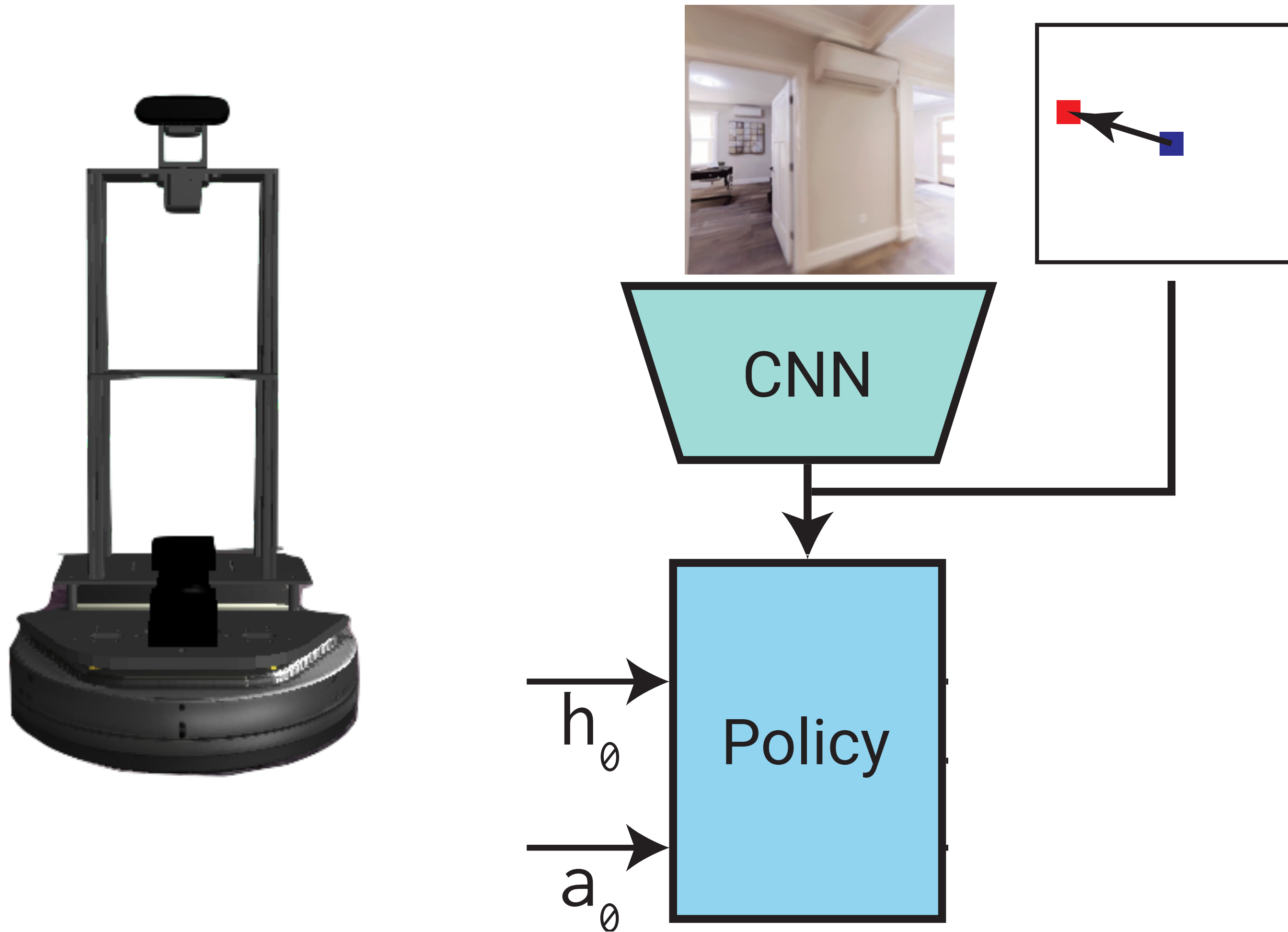
Agent and Model Design



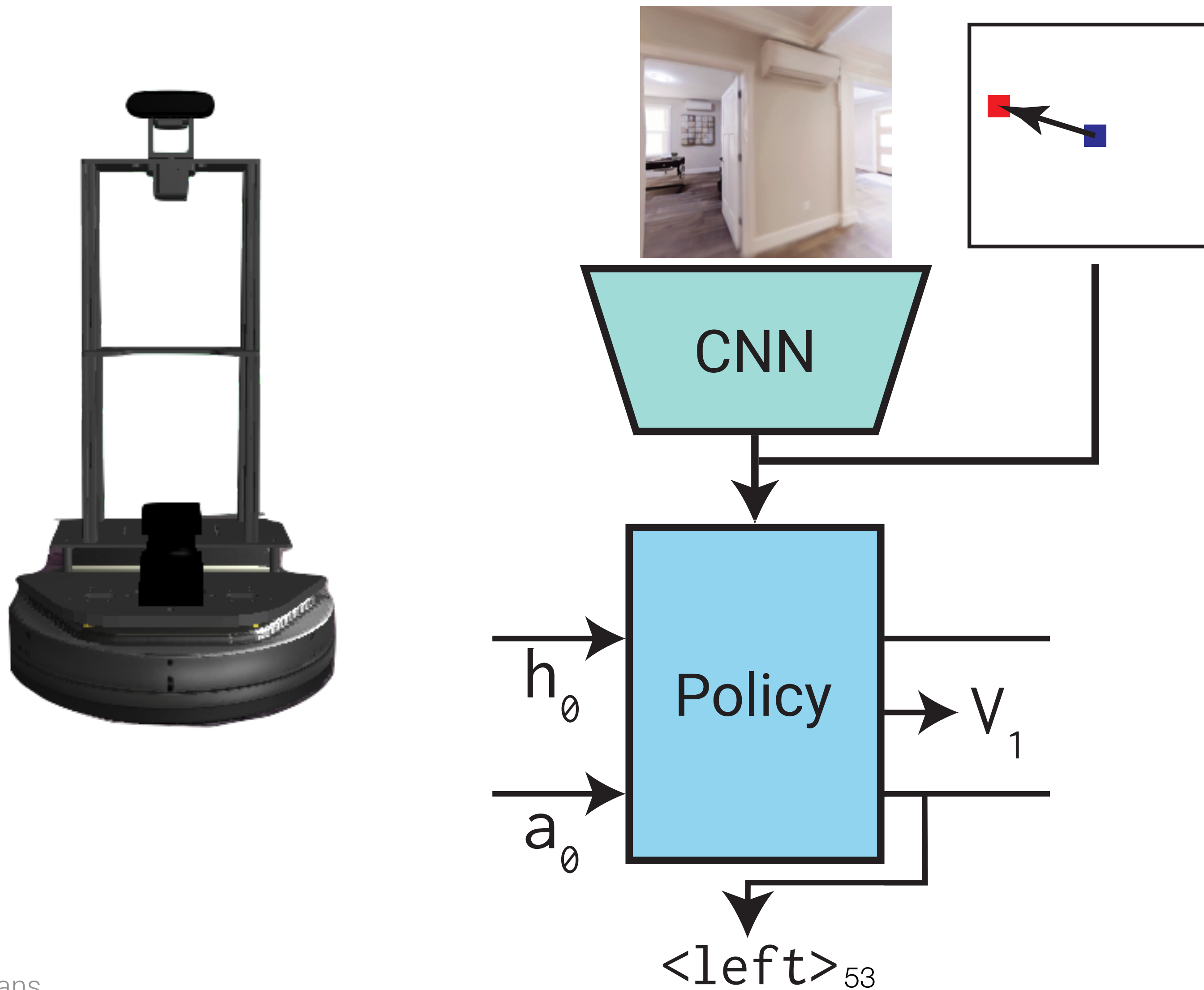
Agent and Model Design



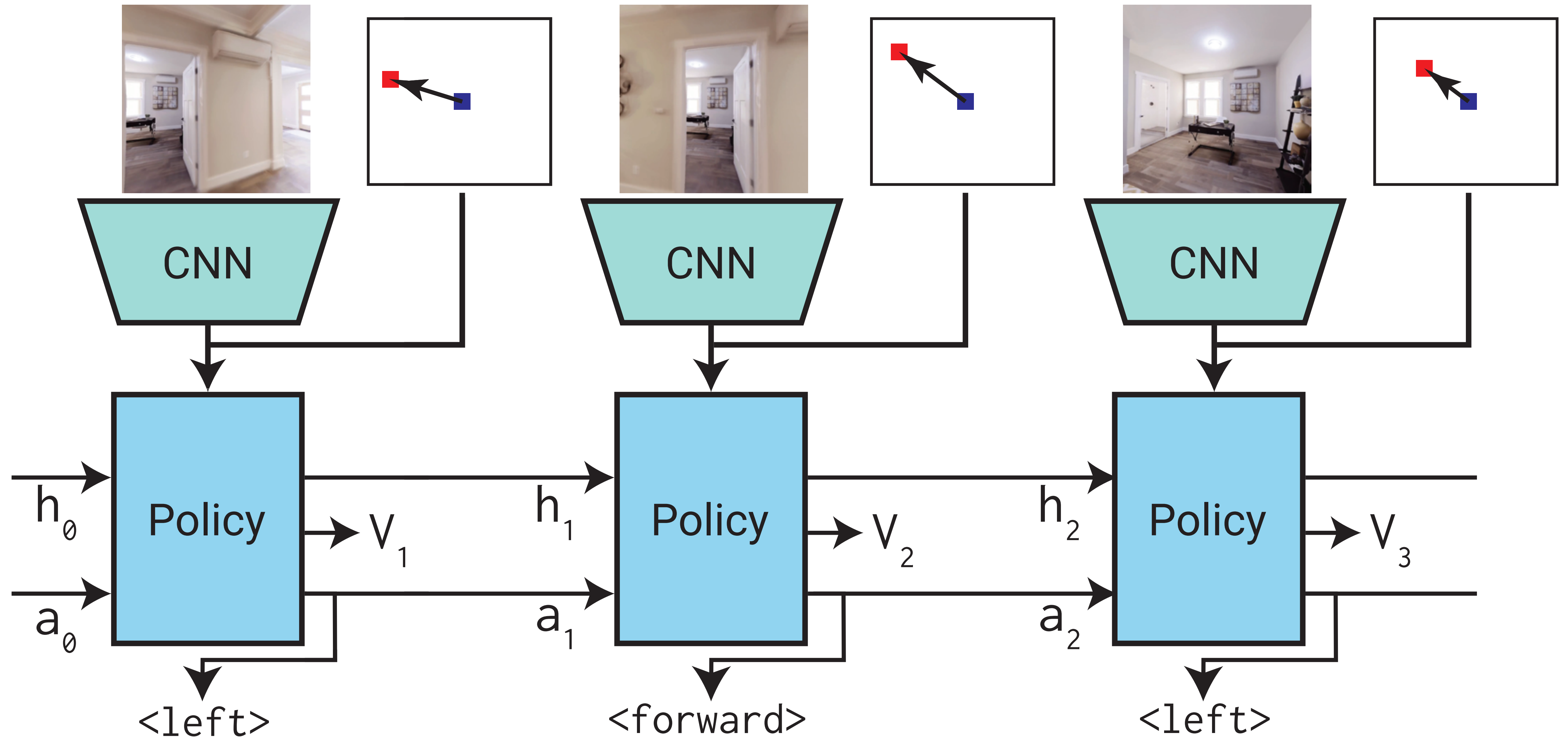
Agent and Model Design



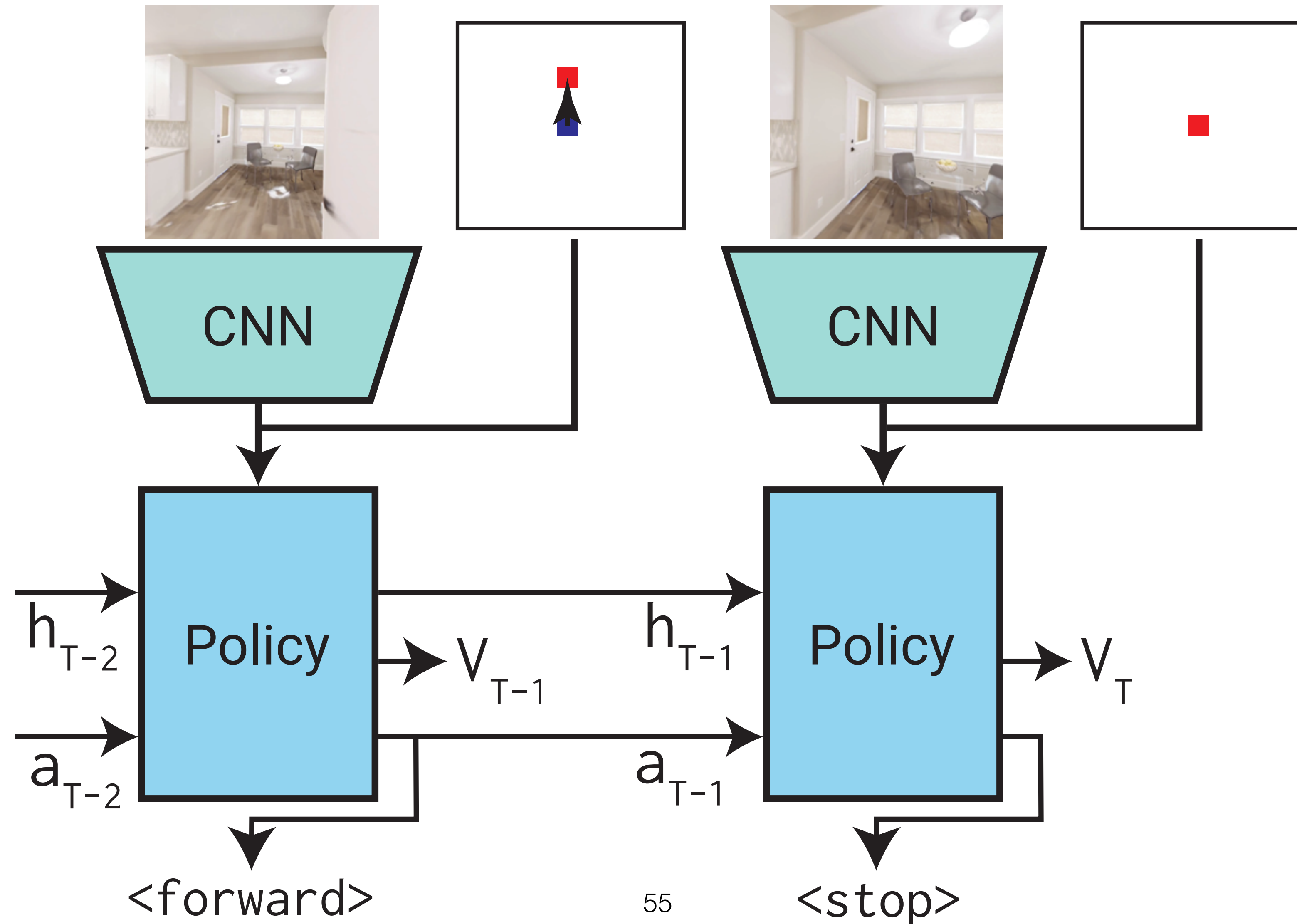
Agent and Model Design



Agent and Model Design



Agent and Model Design



Agent and Model Design



- How do we train this agent?

Agent and Model Design



- How do we train this agent?
- Both actions (they are discrete) and the simulation are non-differential-able

Agent and Model Design



- How do we train this agent?
- Both actions (they are discrete) and the simulation are non-differential-able
- Use reinforcement learning!

Outline

- Proximal Policy Optimization (PPO)
- Application: PointGoal Navigation
 - Sim2Real Transfer
 - Robot2Robot Transfer

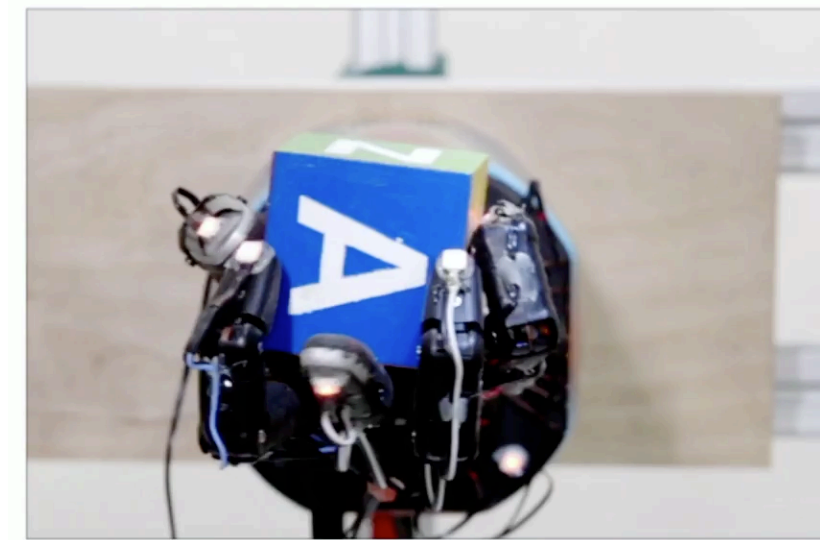
Outline

- Proximal Policy Optimization (PPO)
- Application: PointGoal Navigation
 - Sim2Real Transfer
 - Robot2Robot Transfer

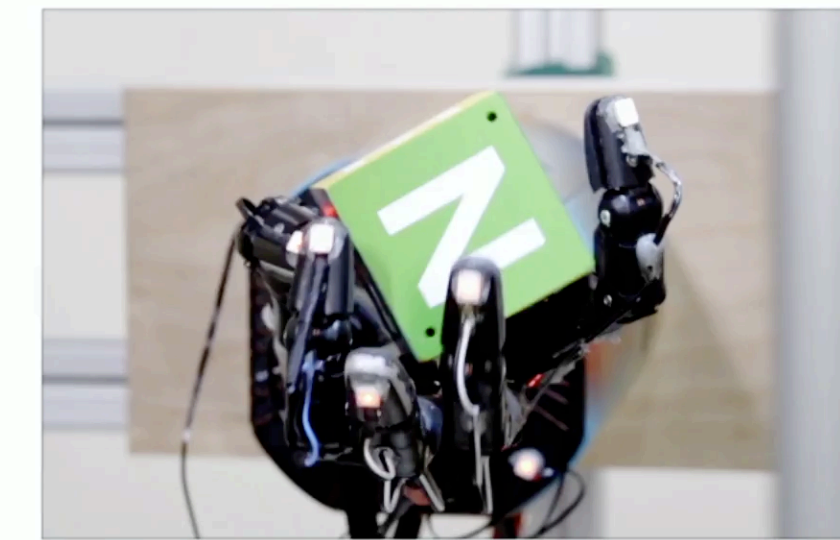
Proximal Policy Optimization (PPO)

Proximal Policy Optimization (PPO)

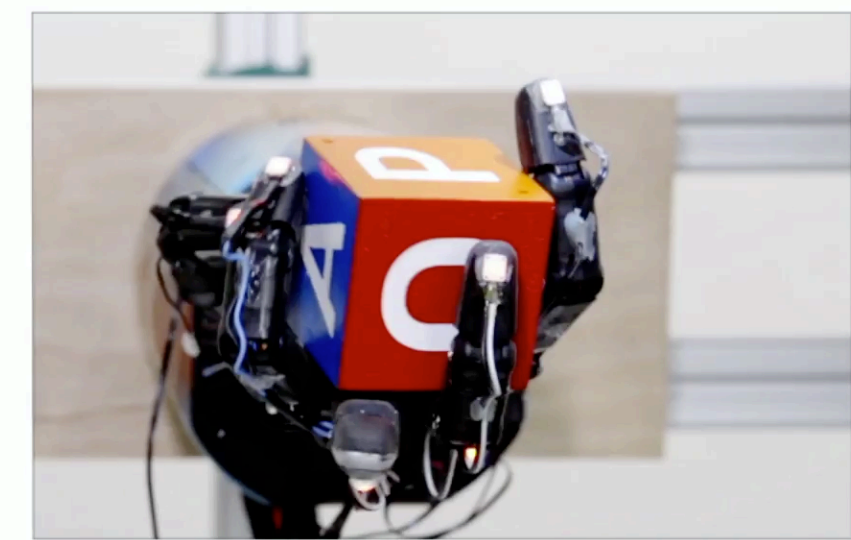
AlphaStar: Mastering the Real-Time Strategy Game StarCraft II



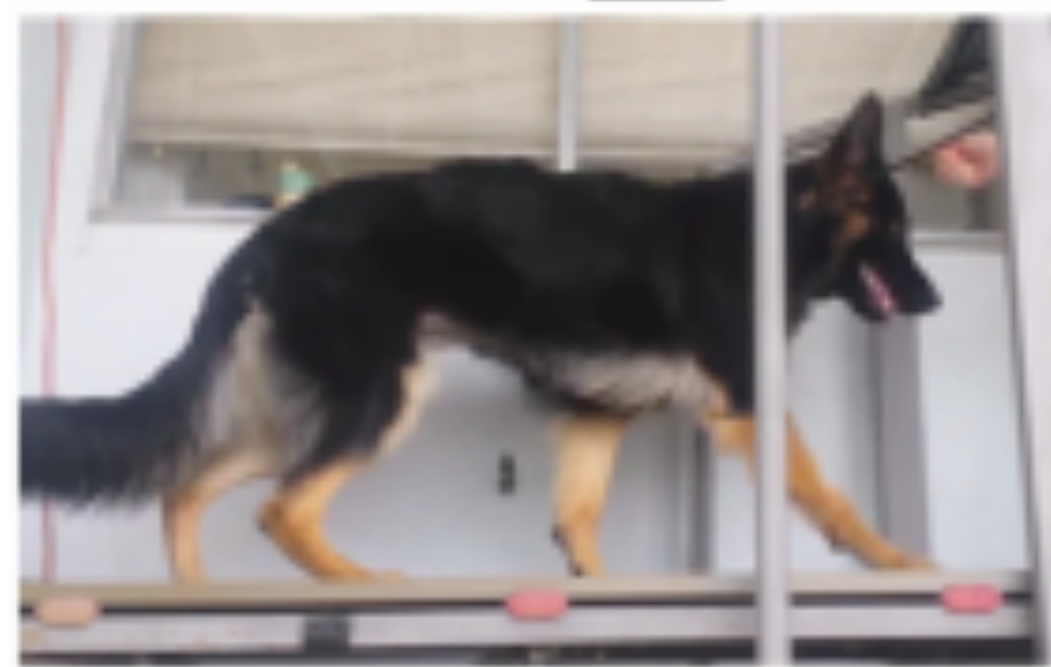
FINGER PIVOTING



SLIDING



FINGER GAITING



Animal



Reference



Simulation



Real Robot

Proximal Policy Optimization (PPO)

Proximal Policy Optimization (PPO)

Given a policy: $\pi_{\theta_{\text{old}}}(a_t | s_t)$

Proximal Policy Optimization (PPO)

Given a policy: $\pi_{\theta_{\text{old}}}(a_t | s_t)$

Objective— maximize probability ratio: $ratio_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$

Proximal Policy Optimization (PPO)

Given a policy: $\pi_{\theta_{old}}(a_t | s_t)$

Objective— maximize probability ratio: $ratio_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}$

Advantage: $A_{\pi_{\theta_{old}}}(s_t, a_t) = Q_{\pi_{\theta_{old}}}(s, a) - V_{\pi_{\theta_{old}}}(s)$

Proximal Policy Optimization (PPO)

Given a policy: $\pi_{\theta_{old}}(a_t | s_t)$

Objective— maximize probability ratio: $ratio_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}$

Advantage: $A_{\pi_{\theta_{old}}}(s_t, a_t) = Q_{\pi_{\theta_{old}}}(s, a) - V_{\pi_{\theta_{old}}}(s)$

$$\mathcal{J}^{PPO}(\theta) = A_{\pi_{\theta_{old}}}(s_t, a_t) \cdot \begin{cases} \min(r_t(\theta), 1 + \epsilon) & \text{if } A_{\pi_{\theta_{old}}}(s_t, a_t) > 0 \\ \max(r_t(\theta), 1 - \epsilon) & \text{if } A_{\pi_{\theta_{old}}}(s_t, a_t) < 0 \end{cases}$$

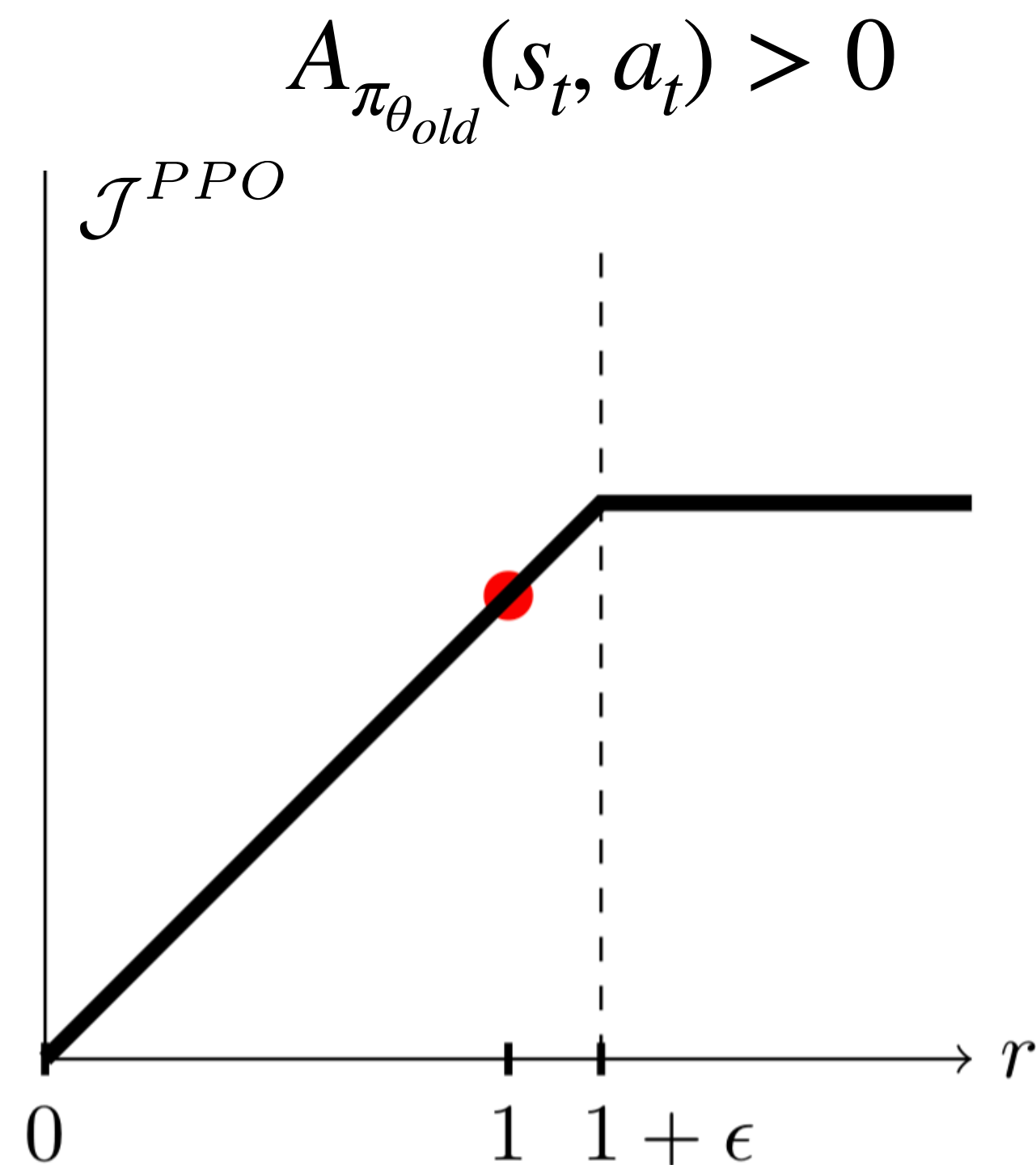
Proximal Policy Optimization (PPO)

$$\mathcal{J}^{\text{PPO}}(\theta) = A_{\pi_{\theta_{old}}}(s_t, a_t) \cdot \begin{cases} \min(r_t(\theta), 1 + \epsilon) & \text{if } A_{\pi_{\theta_{old}}}(s_t, a_t) > 0 \\ \max(r_t(\theta), 1 - \epsilon) & \text{if } A_{\pi_{\theta_{old}}}(s_t, a_t) < 0 \end{cases}$$

Proximal Policy Optimization (PPO)

$$\mathcal{J}^{\text{PPO}}(\theta) = A_{\pi_{\theta_{\text{old}}}}(s_t, a_t) \cdot \begin{cases} \min(r_t(\theta), 1 + \epsilon) & \text{if } A_{\pi_{\theta_{\text{old}}}}(s_t, a_t) > 0 \\ \max(r_t(\theta), 1 - \epsilon) & \text{if } A_{\pi_{\theta_{\text{old}}}}(s_t, a_t) < 0 \end{cases}$$

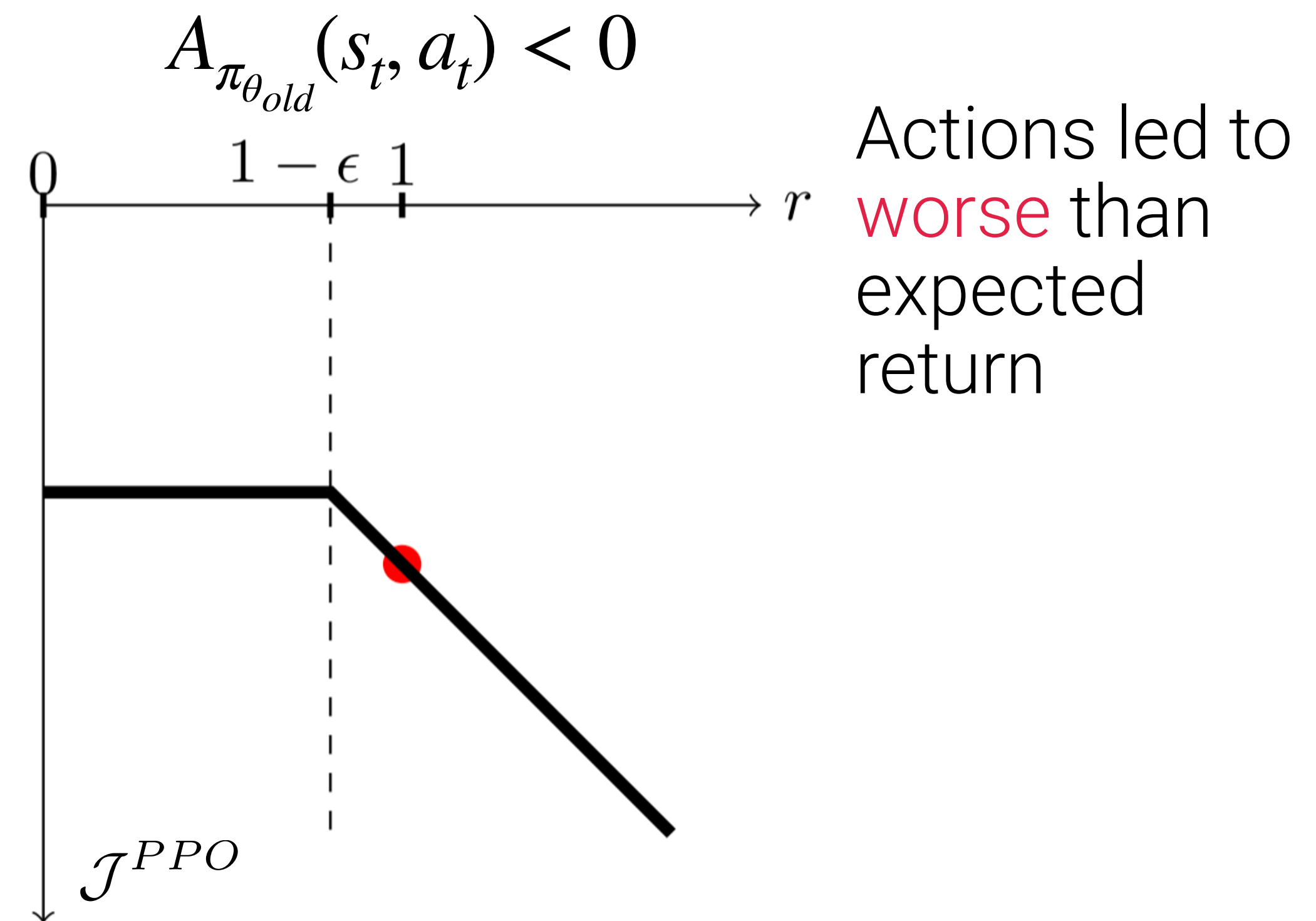
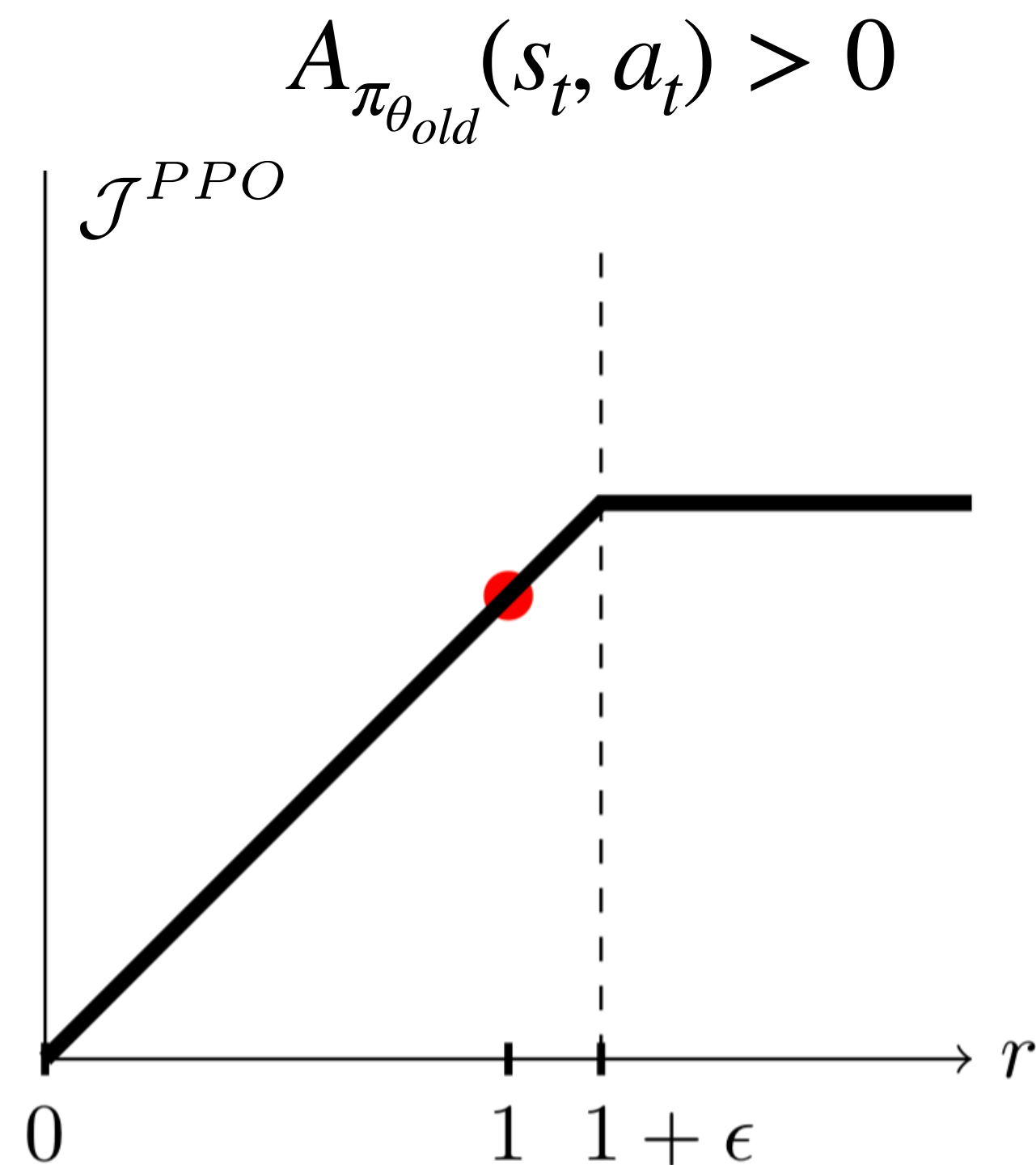
Actions led to
better than
expected
return



Proximal Policy Optimization (PPO)

$$\mathcal{J}^{\text{PPO}}(\theta) = A_{\pi_{\theta_{\text{old}}}}(s_t, a_t) \cdot \begin{cases} \min(r_t(\theta), 1 + \epsilon) & \text{if } A_{\pi_{\theta_{\text{old}}}}(s_t, a_t) > 0 \\ \max(r_t(\theta), 1 - \epsilon) & \text{if } A_{\pi_{\theta_{\text{old}}}}(s_t, a_t) < 0 \end{cases}$$

Actions led to **better** than expected return



Actions led to **worse** than expected return

Proximal Policy Optimization (PPO)

- Advantage
 - Able to perform multiple optimization steps per rollout
 - $\epsilon=0.2$ “just works” in a lot of cases
 - Easily handles networks with hundreds of millions of parameters

Proximal Policy Optimization (PPO)

- Advantage
 - Able to perform multiple optimization steps per rollout
 - $\epsilon=0.2$ “just works” in a lot of cases
 - Easily handles networks with hundreds of millions of parameters
- Disadvantage
 - Other methods are more sample efficient

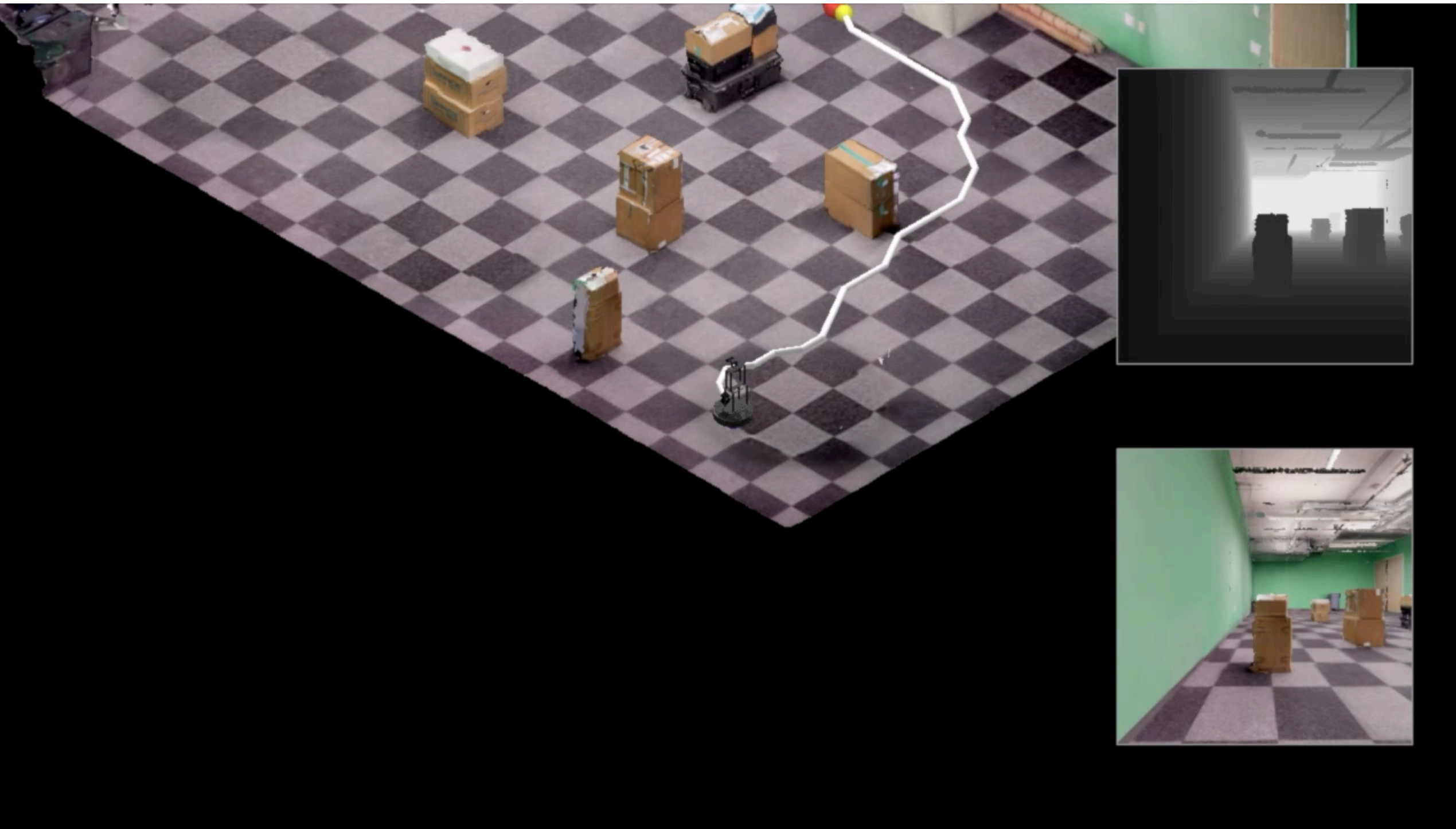
PPO Implementation

1. Collect a set of trajectories using current policy
2. For a few epochs (typically 2 or 4)
 1. Sample mini batches from rollout (typically 2 or 4)
 1. Update the policy via step of PPO/TRPO objective
3. Repeat

Outline

- Proximal Policy Optimization (PPO)
- Soft Actor Critic (SAC)
- Application: PointGoal Navigation
 - Sim2Real Transfer
 - Robot2Robot Transfer

Sim2Real Transfer



Simulation



Reality



●
Goal

20x

Outline

- Proximal Policy Optimization (PPO)
- Soft Actor Critic (SAC)
- Application: PointGoal Navigation
 - Sim2Real Transfer
 - Robot2Robot Transfer

Dynamics Aware Navigation

Dynamics Aware Navigation



A1



AlienGo



Daisy

How to generalize to new robots?



Laikago



Daisy

Sphere baseline (no dynamics)

Sphere baseline (no dynamics)

- Idealized agent: sphere
- Given direct access to points along the shortest path to the goal

Sphere baseline (no dynamics)

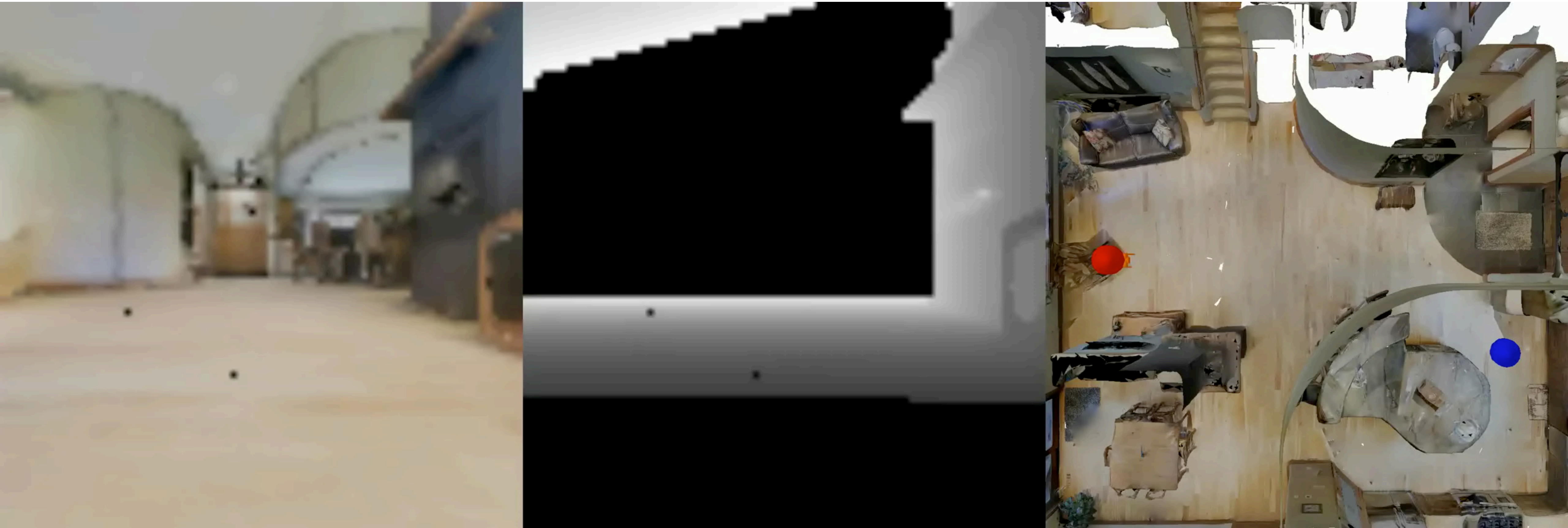


Zero shot transfer: Daisy \rightarrow AlienGo

Zero shot transfer: Daisy \rightarrow AlienGo

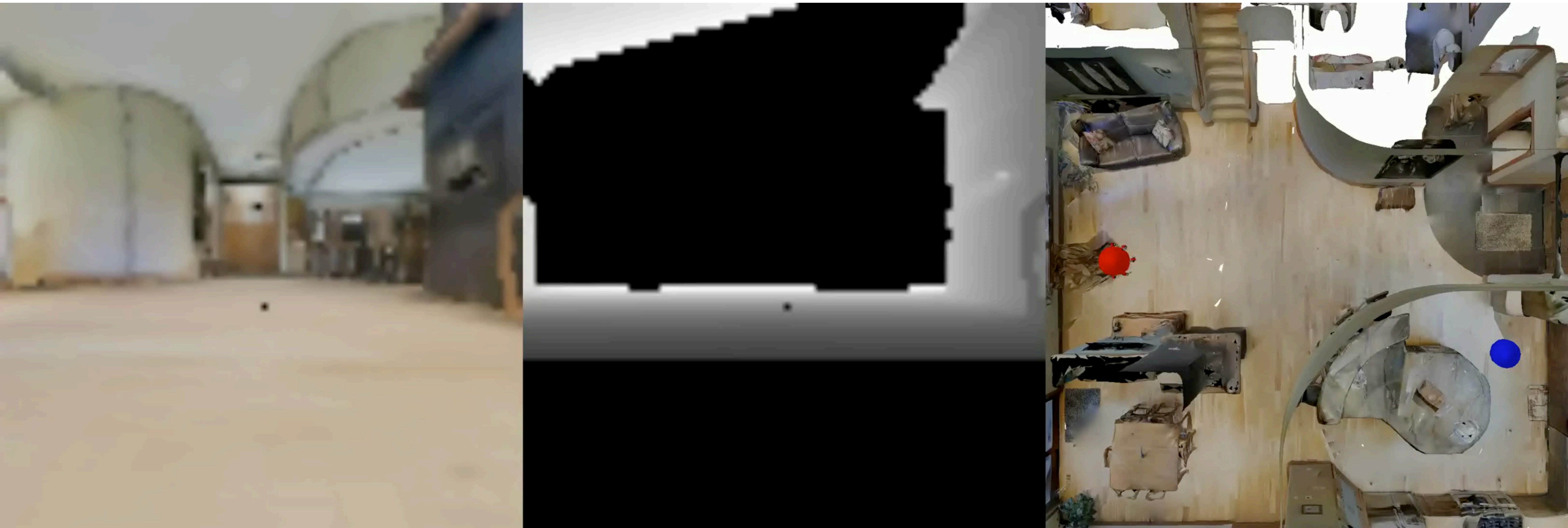
- Trained policy on Daisy robot, deploy on AlienGo

Zero shot transfer: Daisy \rightarrow AlienGo

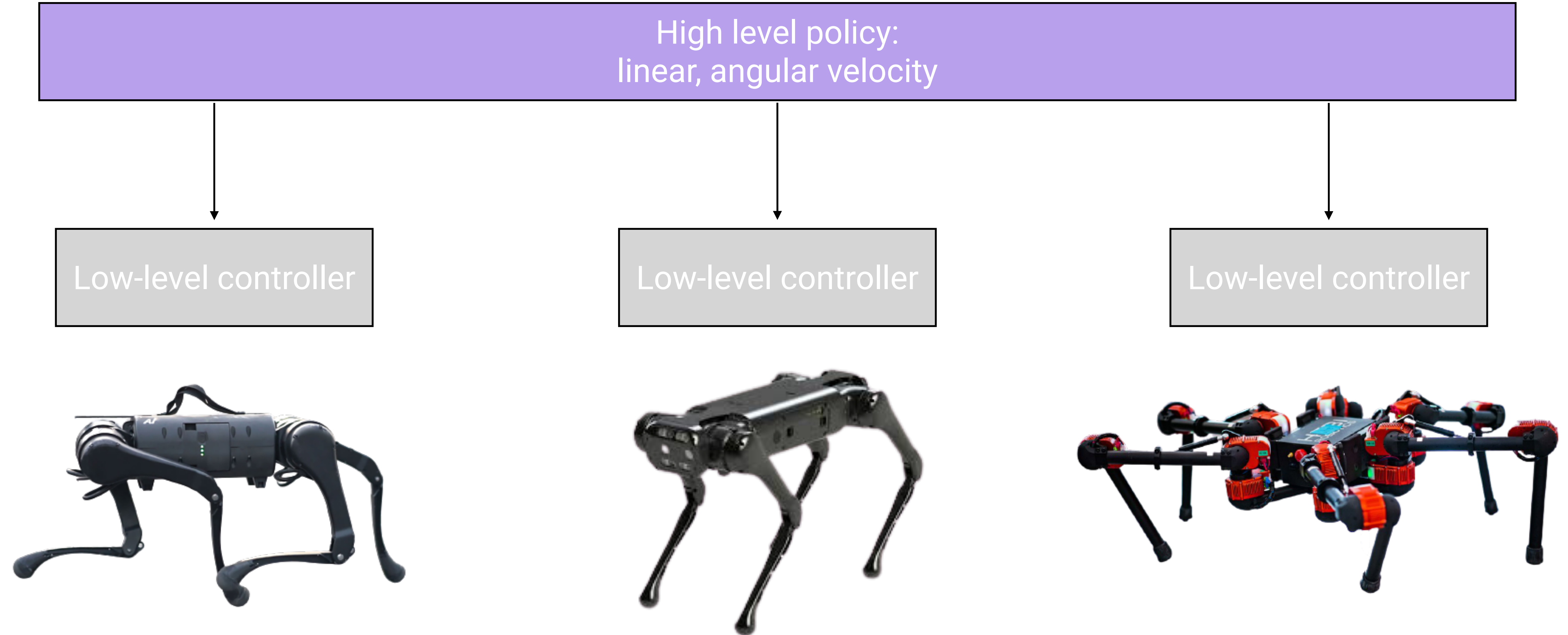


Dynamics aware navigation

Dynamics aware navigation

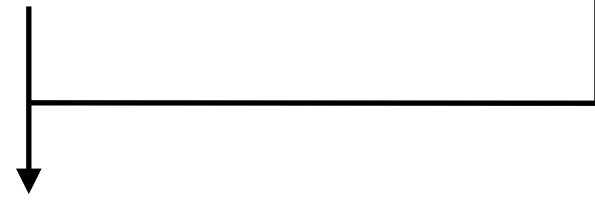
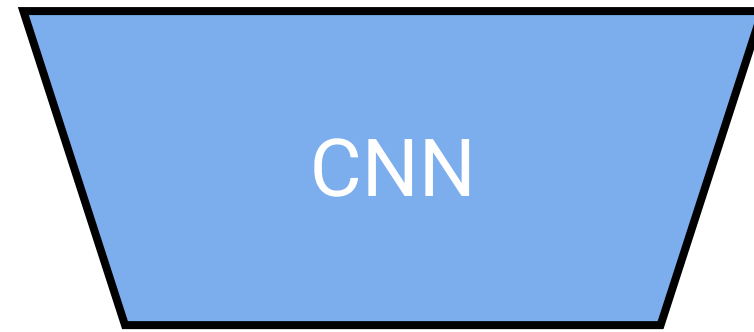
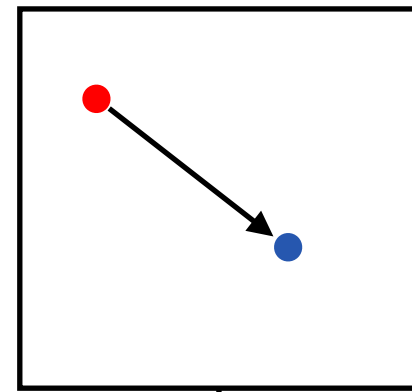


Hierarchical Reinforcement Learning

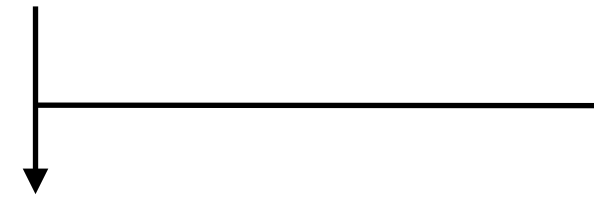
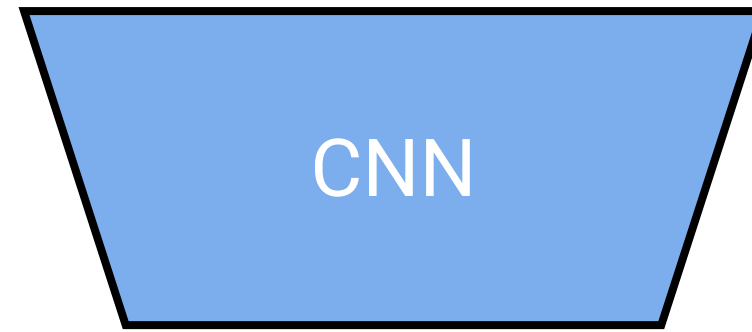
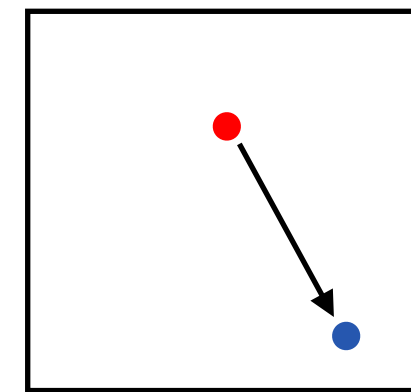


Our approach

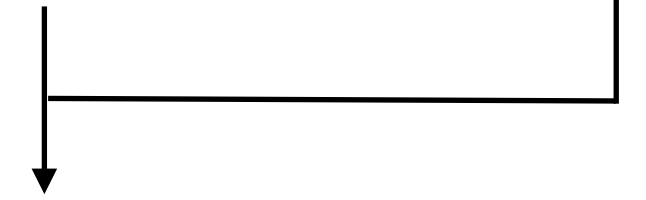
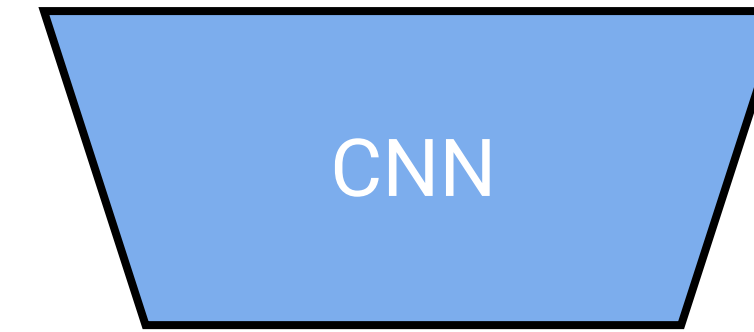
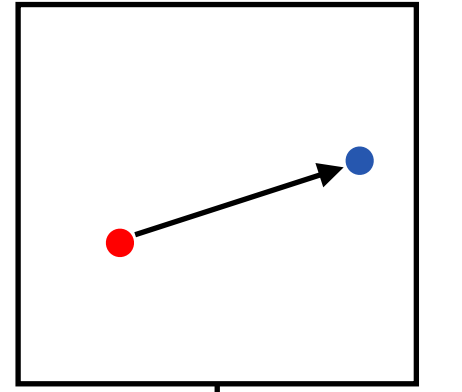
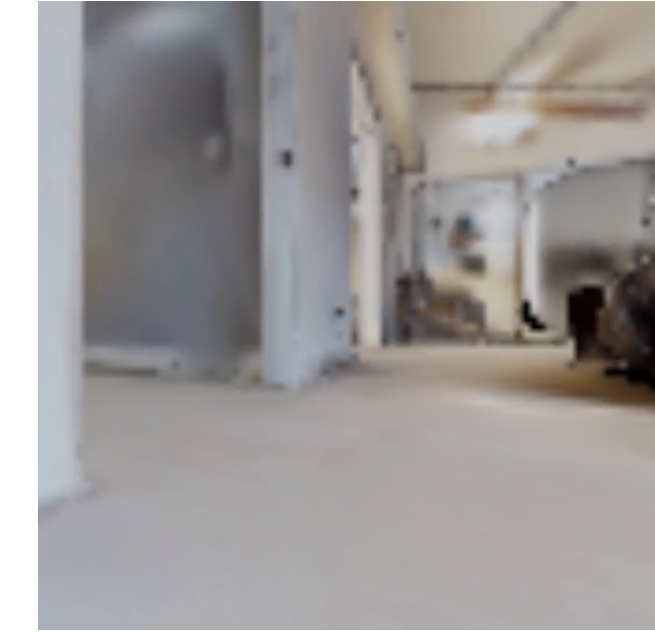
A1



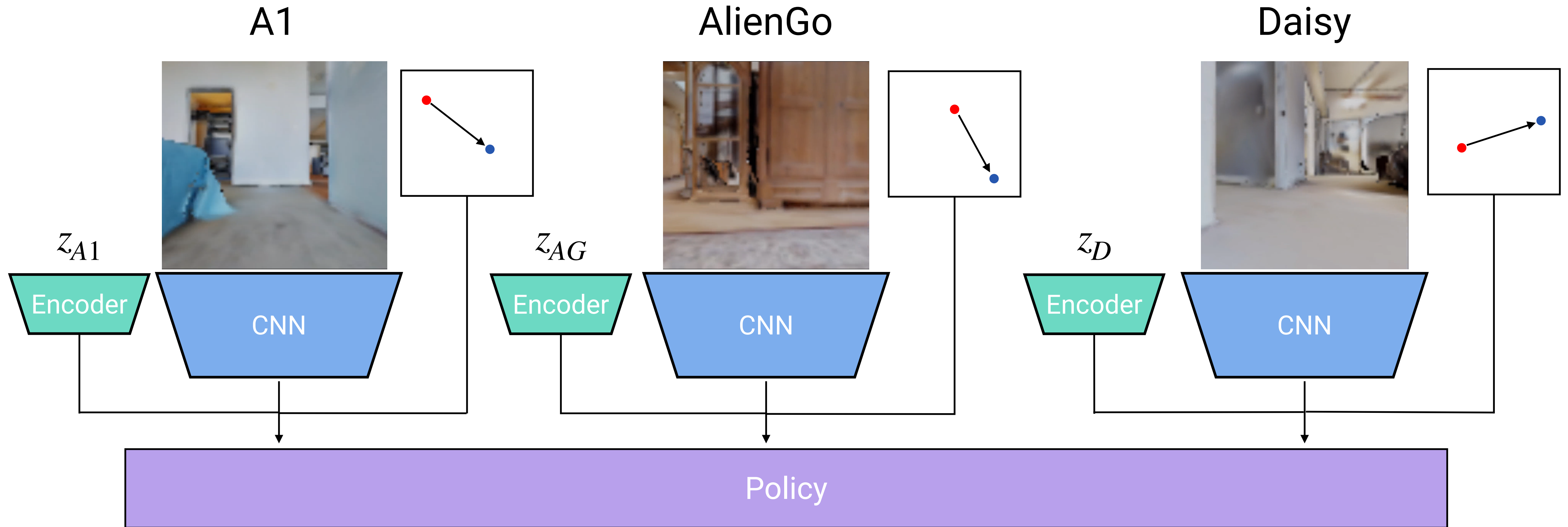
AlienGo



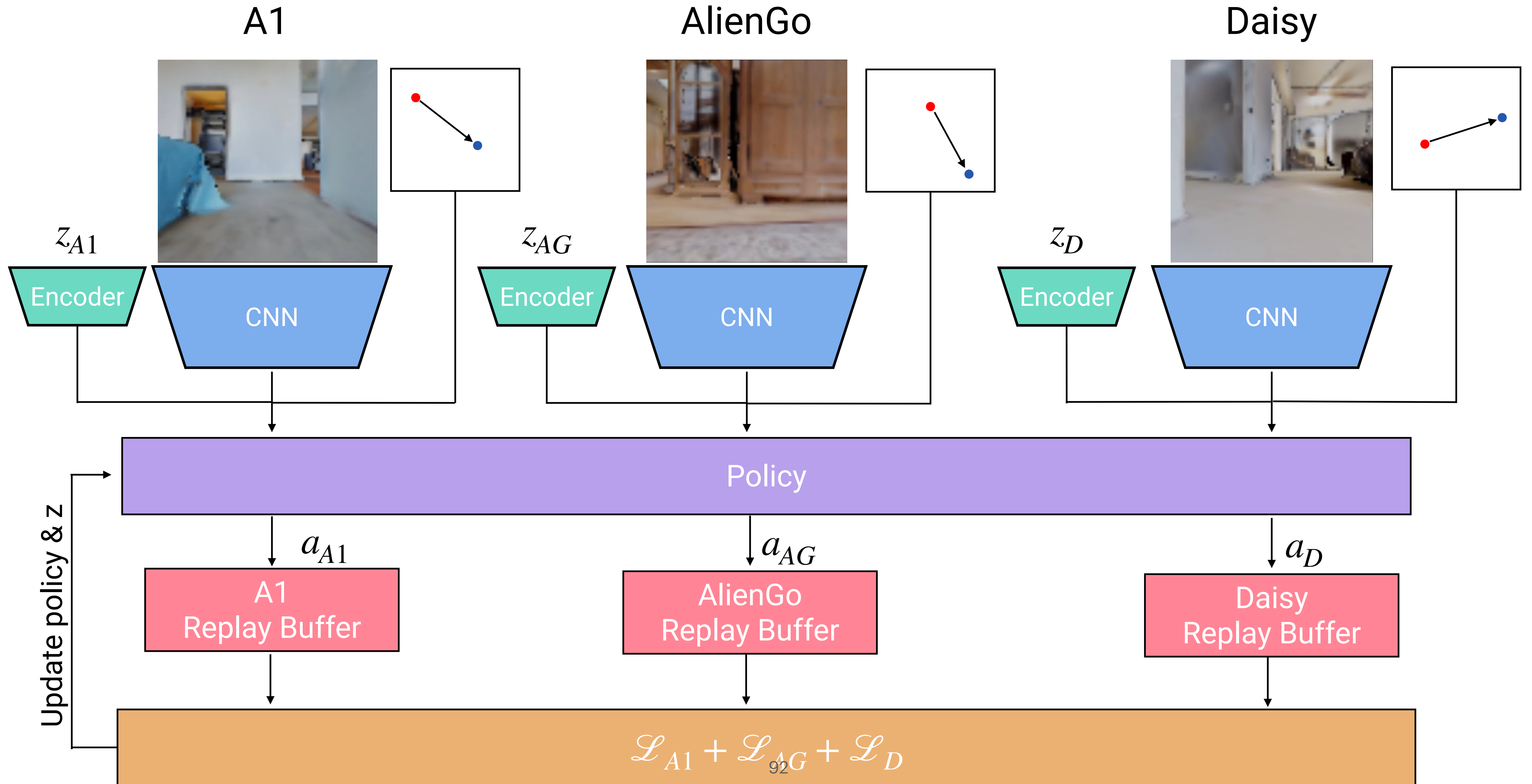
Daisy



Our approach

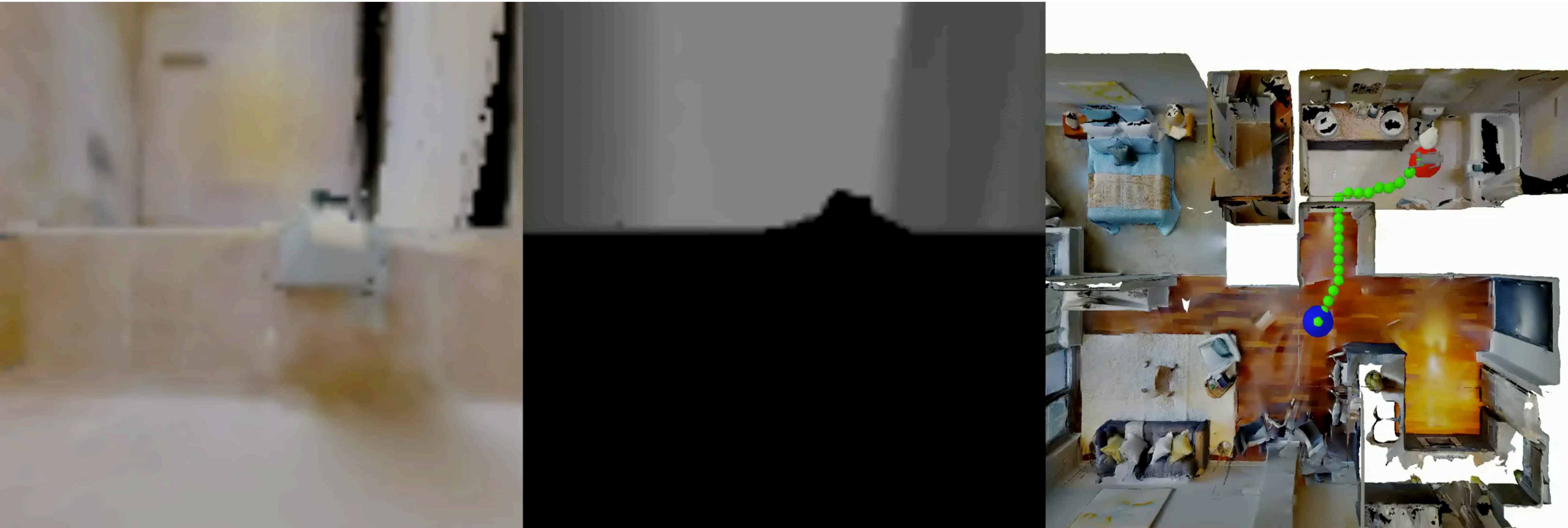


Our approach



A1 (Train robot) in novel environment

A1 (Train robot) in novel environment



AlienGo (Train robot) in novel environment

AlienGo (Train robot) in novel environment



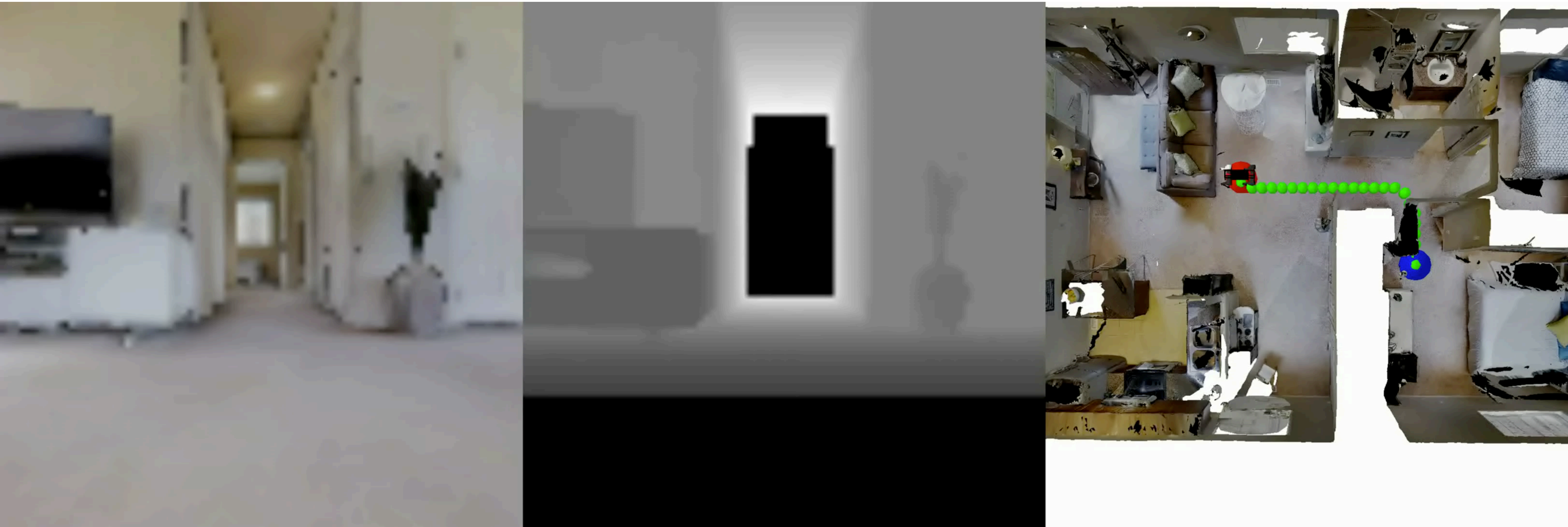
Daisy (Train robot) in novel environment

Daisy (Train robot) in novel environment



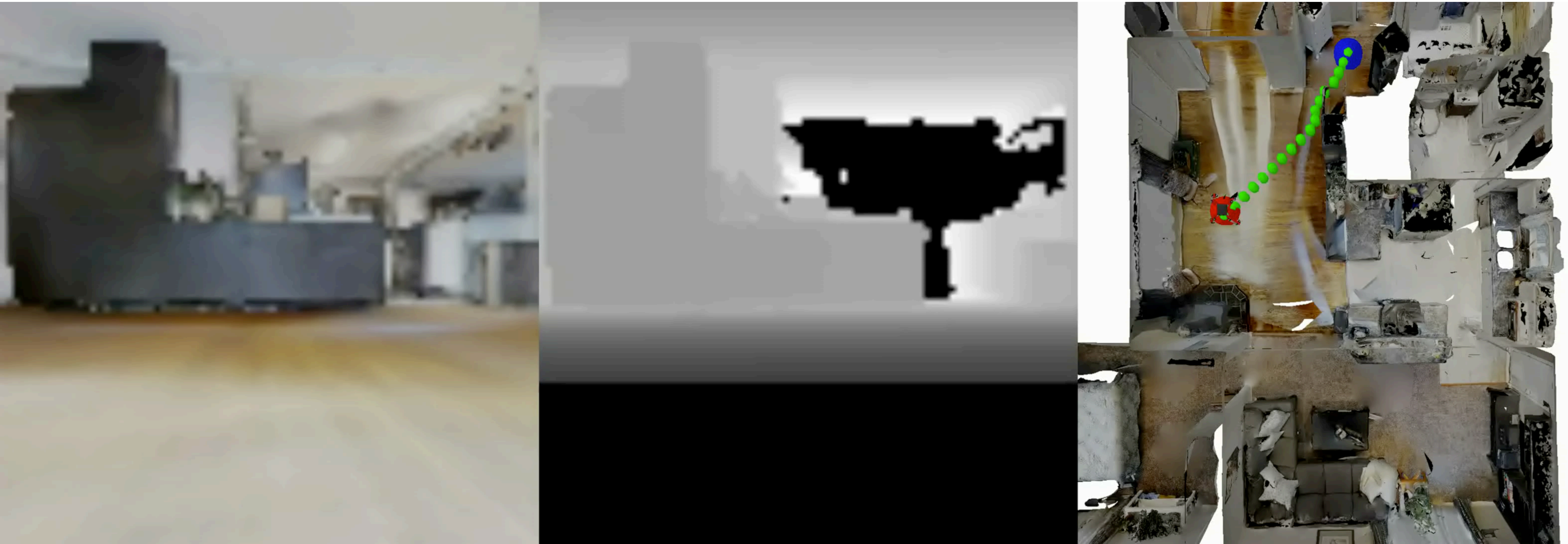
Laikago (new robot) in novel environment

Laikago (new robot) in novel environment



Daisy-4Legged (new robot) in novel environment

Daisy-4Legged (new robot) in novel environment



Thank you!

Questions?