# Last-Mile Embodied Visual Navigation

In submission to CoRL 2022

EAI seminar | Aug 19th 2022

**Justin Wasserman***  Karmesh Yadav  Girish Chowdhary  Abhinav Gupta  Unnat Jain*

UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN  Carnegie Mellon University  Meta AI

# Image-goal navigation

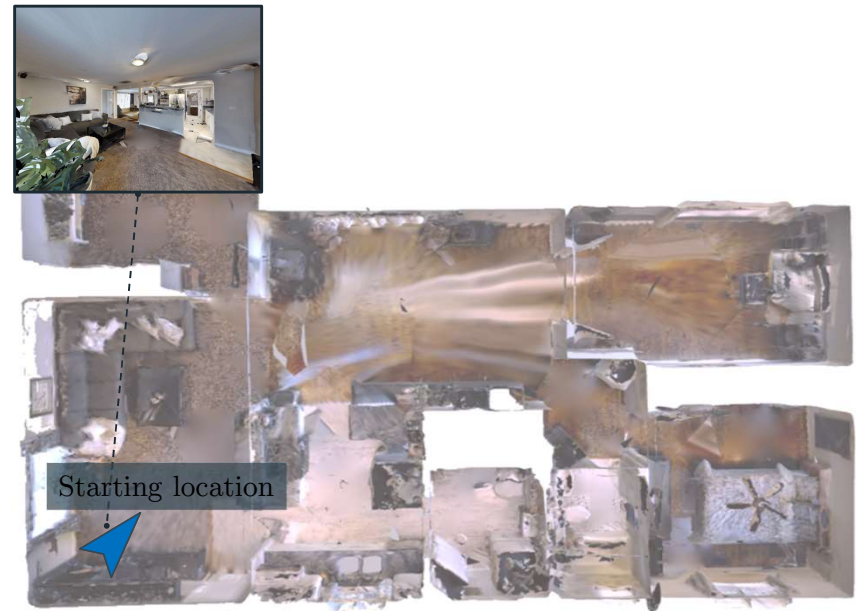❖Observes RGBD and pose at each step



Starting location

# Image-goal navigation

❖Observes RGBD and pose at each step



Starting location

# Image-goal navigation

❖Observes RGBD and pose at each step

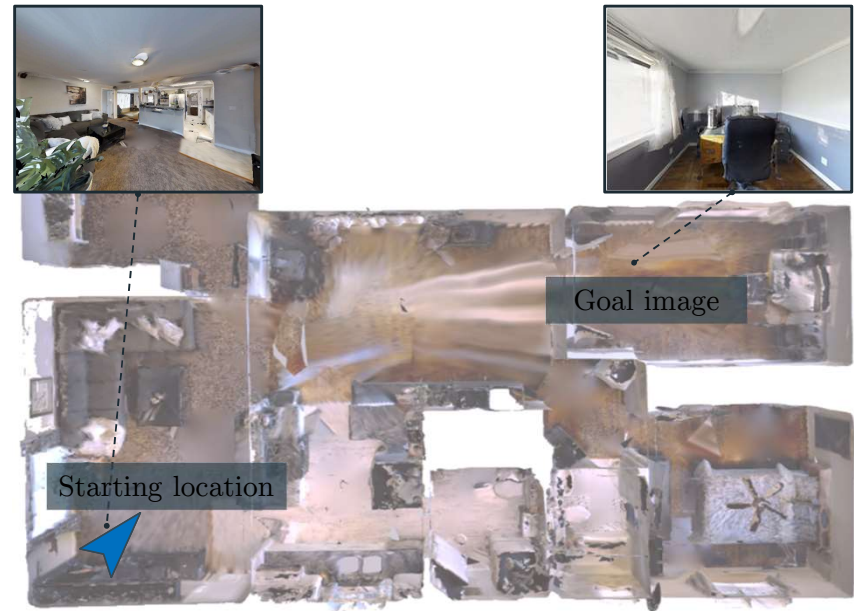❖Must navigate to goal RGB image



Goal image

Starting location

# Image-goal navigation

❖Observes RGBD and pose at each step

❖Must navigate to goal RGB image

❖Action space: {forward, turn right, turn left, stop}



Goal image

Starting location

# Image-goal navigation

❖ Observes RGBD and pose at each step

❖ Must navigate to goal RGB image
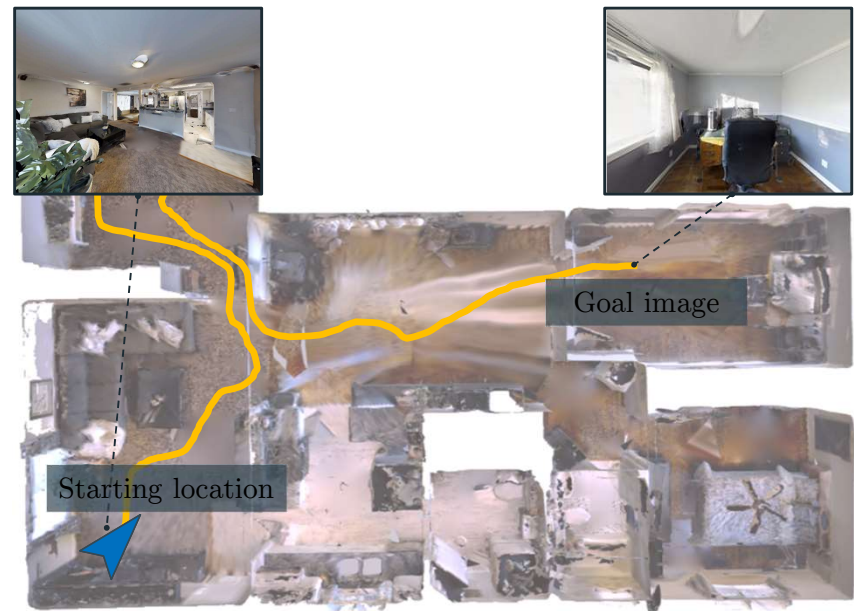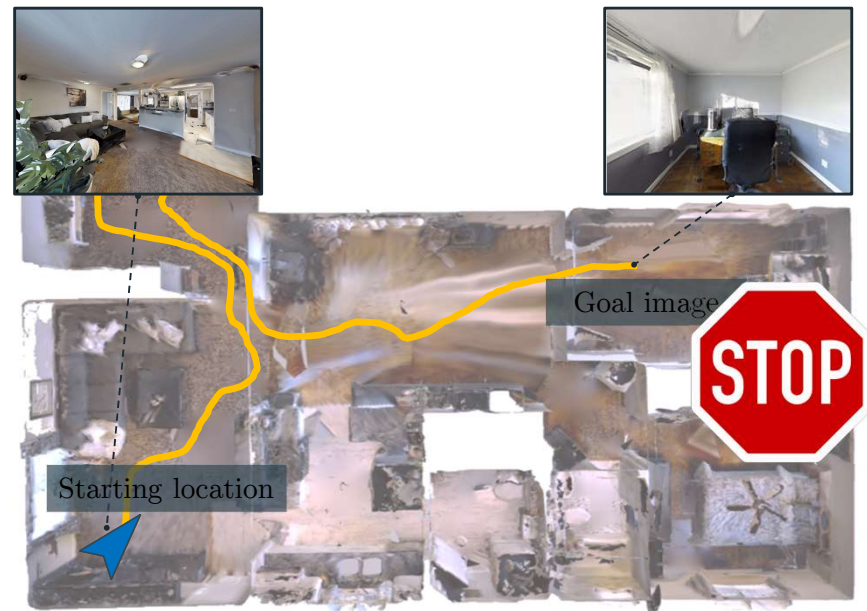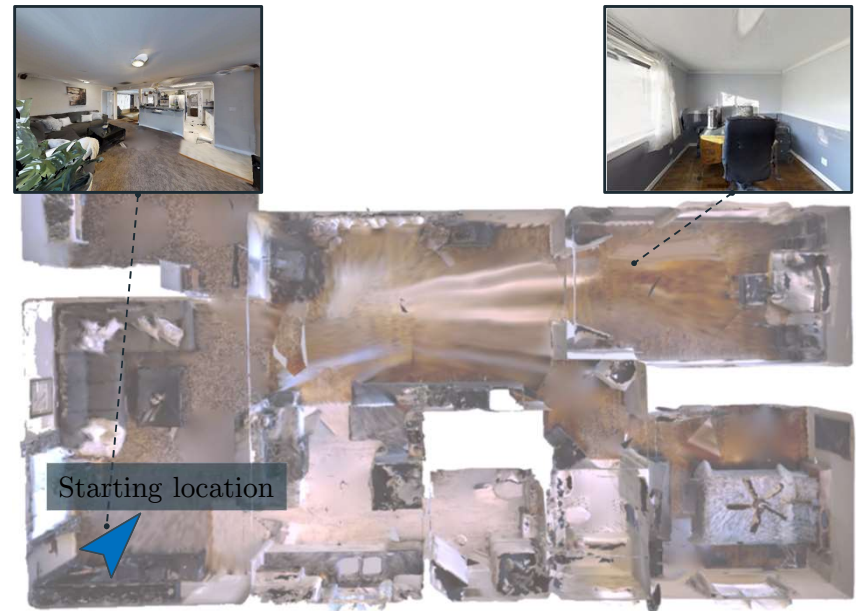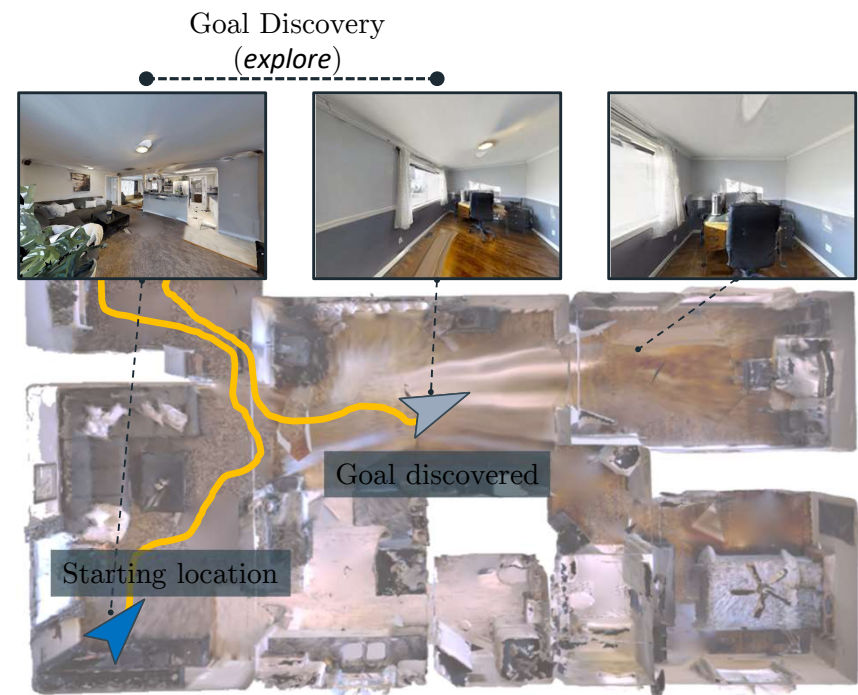
❖ Action space: {forward, turn right, turn left, stop}
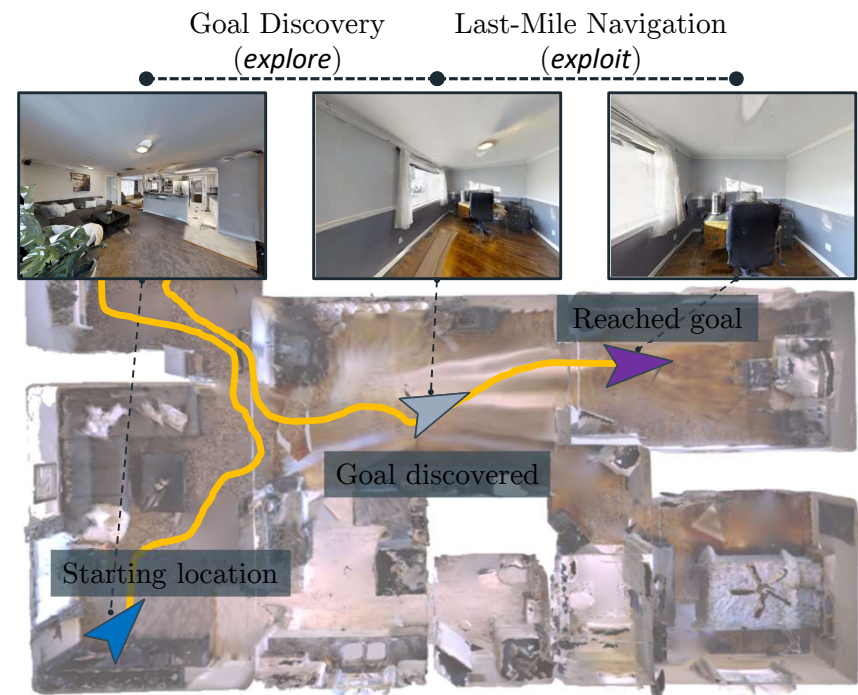
❖ Success: "stop" within 1m of the goal within 500 steps
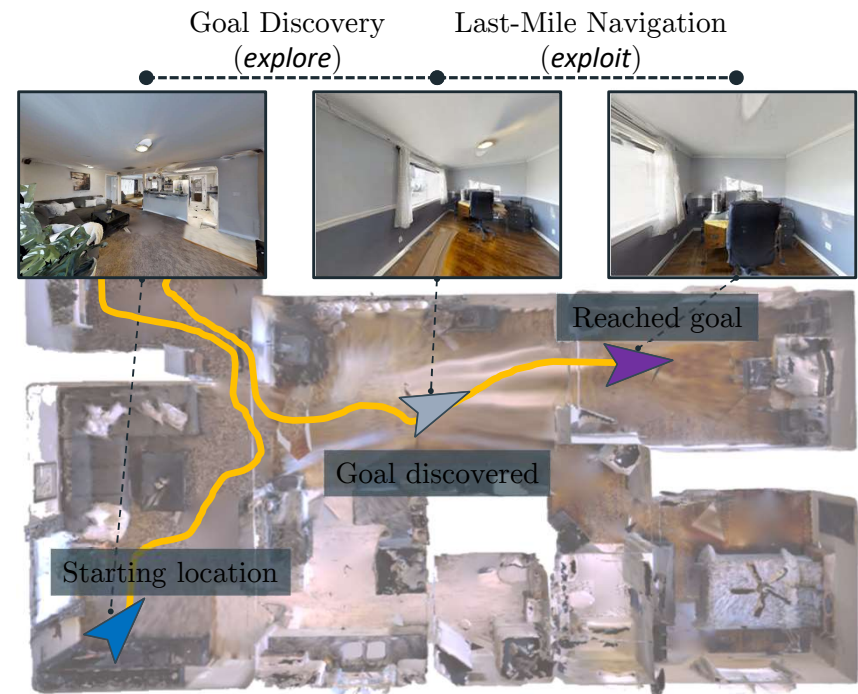
# Motivation

# Motivation



Goal Discovery
(*explore*)

Goal discovered

Starting location

# Motivation

# Motivation

Agents across embodied AI have difficulties getting close to the goal even when they see the goal.



Goal Discovery (*explore*)  Last-Mile Navigation (*exploit*)

Reached goal

Goal discovered

Starting location

# Motivation

S. Wani et al. **MultiON**
[NeurIPS 2020]

FOUND terminates the episode immediately, we allow the agent to call a fixed number of wrong FOUND actions during the episode. We found that allowing even a single wrong FOUND action leads to a significant increase in performance metrics. This suggests that many episodes terminate due to calling FOUND action at the wrong time and fixing this inadequacy could improve $m$-ON performance significantly. Table 9 summarizes the results of performance metrics against the number of wrong FOUND actions allowed in OracleMap model on the 3-ON task. Note that these evaluations were

# Motivation

S. Wani et al. **MultiON**
[NeurIPS 2020]

P. Chattopadhyay et al. **RobustNav**
[ICCV 2021]

FOUND terminates the episode immediately, we allow the agent to call a fixed number of wrong FOUND actions during the episode. We found that allowing even a single wrong FOUND action leads to a significant increase in performance metrics. This suggests that many episodes terminate due to calling FOUND action at the wrong time and fixing this inadequacy could improve $m$-ON performance significantly. Table 9 summarizes the results of performance metrics against the number of wrong FOUND actions allowed in OracleMap model on the 3-ON task. Note that these evaluations were
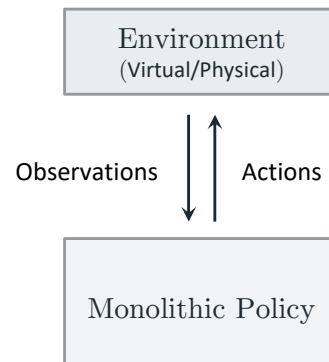
| Last mile | Gets stuck near the goal. |
| Commitment | Sees and approaches the goal but passes it. |
| Open | Explores an open area without any objects. |

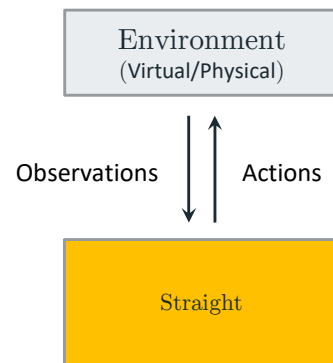**Table 4.** Description of prominent failure modes.

# Motivation

S. Wani et al. **MultiON**
[NeurIPS 2020]

FOUND terminates the episode immediately, we allow the agent to call a fixed number of wrong FOUND actions during the episode. We found that allowing even a single wrong FOUND action leads to a significant increase in performance metrics. This suggests that many episodes terminate due to calling FOUND action at the wrong time and fixing this inadequacy could improve $m$-ON performance significantly. Table 9 summarizes the results of performance metrics against the number of wrong FOUND actions allowed in OracleMap model on the 3-ON task. Note that these evaluations were

P. Chattopadhyay et al. **RobustNav**
[ICCV 2021]

| Last mile | Gets stuck near the goal. |
| Commitment | Sees and approaches the goal but passes it. |
| Open | Explores an open area without any objects. |

**Table 4.** Description of prominent failure modes.

J. Ye et al. **ObjectAux**
[ICCV 2021]

**Corruptions hurt OBJECTNAV stopping mechanism.** Recall that for both POINTNAV and OBJECTNAV, success depends on the notion of "intentionality" [5] – the agent calls an **end** action when it believes it has reached the goal. In Fig 4 (*last two columns*) we aim to understand how cor-
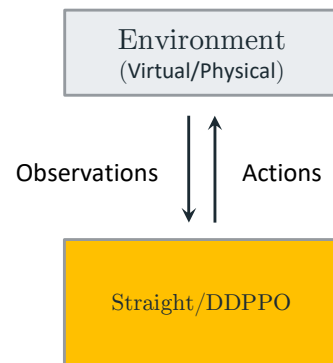
# Solving image-goal navigation

# Solving image-goal navigation

[Heuristic]
T. Chen et al. Learning exploration policies for navigation.
ICLR 2019

```
┌─────────────────────────┐
│      Environment         │
│    (Virtual/Physical)    │
└─────────────────────────┘
           │    ↑
Observations│    │ Actions
           ↓    │
┌─────────────────────────┐
│        Straight          │
│                          │
└─────────────────────────┘
```
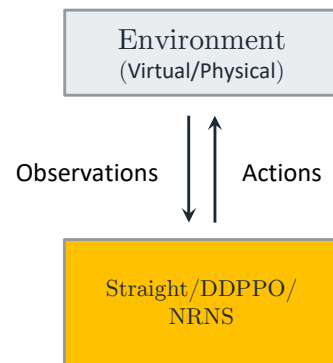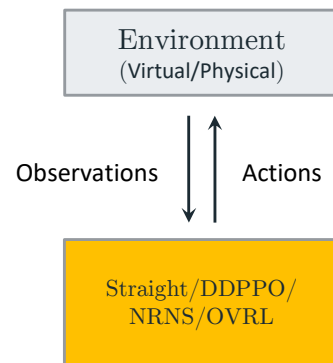
# Solving image-goal navigation

Environment
(Virtual/Physical)

Observations    Actions

Straight/DDPPO

[Heuristic]
T. Chen et al. Learning exploration policies for navigation.
ICLR 2019

[End-to-end deep RL]
E. Wijmans et al. DDPPO
ICLR 2019

# Solving image-goal navigation

Environment
(Virtual/Physical)

Observations    Actions

Straight/DDPPO/
NRNS

[Heuristic]
T. Chen et al. Learning exploration policies for navigation.
ICLR 2019

[End-to-end deep RL]
E. Wijmans et al. DDPPO
ICLR 2019

[Modular]
M. Hahn et al. No rl, no simulation
NeurIPS 2021

# Solving image-goal navigation

Environment
(Virtual/Physical)

Observations          Actions

Straight/DDPPO/
NRNS/OVRL

[Heuristic]
T. Chen et al. Learning exploration policies for navigation.
ICLR 2019

[End-to-end deep RL]
E. Wijmans et al. DDPPO
ICLR 2019

[Modular]
M. Hahn et al. No rl, no simulation
NeurIPS 2021

[SSL]
K. Yadav et al. Offline visual repr. Learning for embodied navigation.
arXiv 2022

# Solving image-goal navigation

Environment
(Virtual/Physical)

Observations     Actions

Straight/DDPPO/
NRNS/OVRL/Oracle

[Heuristic]
T. Chen et al. Learning exploration policies for navigation.
ICLR 2019

[End-to-end deep RL]
E. Wijmans et al. DDPPO
ICLR 2019

[Modular]
M. Hahn et al. No rl, no simulation
NeurIPS 2021

[SSL]
K. Yadav et al. Offline visual repr. Learning for embodied navigation.
arXiv 2022

# Focus on Last-Mile navigation

# Focus on Last-Mile navigation



Environment
(Virtual/Physical)

Observations | Actions

Goal Discovery

Straight/DDPPO/
NRNS/OVRL/Oracle

Previous works

Last-Mile
Navigation

Principled 3D vision module for the last mile

# Focus on Last-Mile navigation



Principled 3D vision module for the last mile

Simple switches to plug-and-play

# Focus on Last-Mile navigation



Principled 3D vision module for the last mile

Simple switches to plug-and-play

# Switchable Last-Mile Image-Goal Navigation (SLING)



Principled 3D vision module for the last mile

Simple switches to plug-and-play

# Switchable Last-Mile Image-Goal Navigation (SLING)



Principled 3D vision module for the last mile

Simple switches to plug-and-play

# Our simple solution to Last-Mile Navigation

# Our simple solution to Last-Mile Navigation

### Agent Image



### Goal Image

# Our simple solution to Last-Mile Navigation



Agent Image ⟶

Goal Image ⟶

# Our simple solution to Last-Mile Navigation

# Our simple solution to Last-Mile Navigation

# Our simple solution to Last-Mile Navigation

# Our simple solution to Last-Mile Navigation

# Our simple solution to Last-Mile Navigation



Feature Extraction · Feature Matcher · Projecting 2D→3D · Perspective-n-Point

Agent Image

Goal Image

Camera Intrinsics

Depth mask

Estimate rotation and translation

M. A. Fischler et al. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography.
ACM 1981

C.-P. Lu et al. Fast and globally convergent pose estimation from video images.
TPAMI 2000.

V. Lepetit et al. An accurate o(n) solution to the pnp problem.
IJCV 2009.

# Our simple solution to Last-Mile Navigation

# Our simple solution to Last-Mile Navigation



Feature Extraction · Feature Matcher · Projecting 2D→3D · Perspective-n-Point

Agent Image · Goal Image · Camera Intrinsics · Depth mask

Estimate rotation and translation

Local Policy

move forward / turn right / turn left / stop

# SOTA on Image-goal navigation in AIHabitat

## SPL

■ w/o SLING

■ w/ SLING

# SOTA on Image-goal navigation in AIHabitat

DDPPO
[End-to-end RL]

SPL

w/o SLING

w/ SLING

Gibson

12.9

22.8

Last-Mile Embodied Visual Navigation

# SOTA on Image-goal navigation in AIHabitat

DDPPO
[End-to-end RL]

NRNS
[Modular IL]

SPL

w/o SLING

w/ SLING

Gibson

12.9    22.8    8.1    15.1

# SOTA on Image-goal navigation in AIHabitat

Last-Mile Embodied Visual Navigation

# SLING transferred to a robot

Last-Mile Embodied Visual Navigation

# SLING transferred to a robot

# SLING transferred to a robot

# SLING transferred to a robot

# SLING transferred to a robot



## SPL

🟧 w/o SLING

🟦 w/ SLING

Real-world, Easy
(1.5 – 3m)

Real-world, Hard
(5 – 10m)
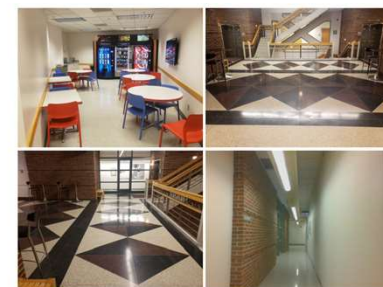
# SLING transferred to a robot



**NRNS**

SPL

- w/o SLING
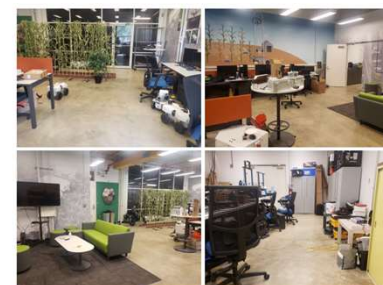- w/ SLING

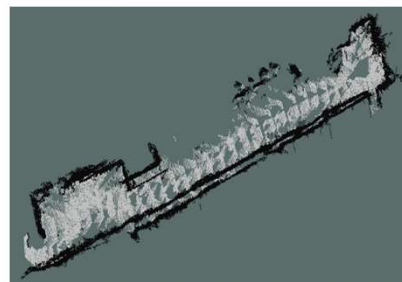37.7   53.7
Real-world, Easy
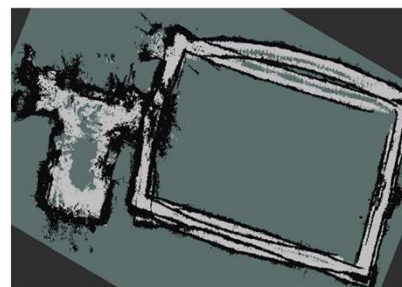(1.5 – 3m)

3.3   19.3
Real-world, Hard
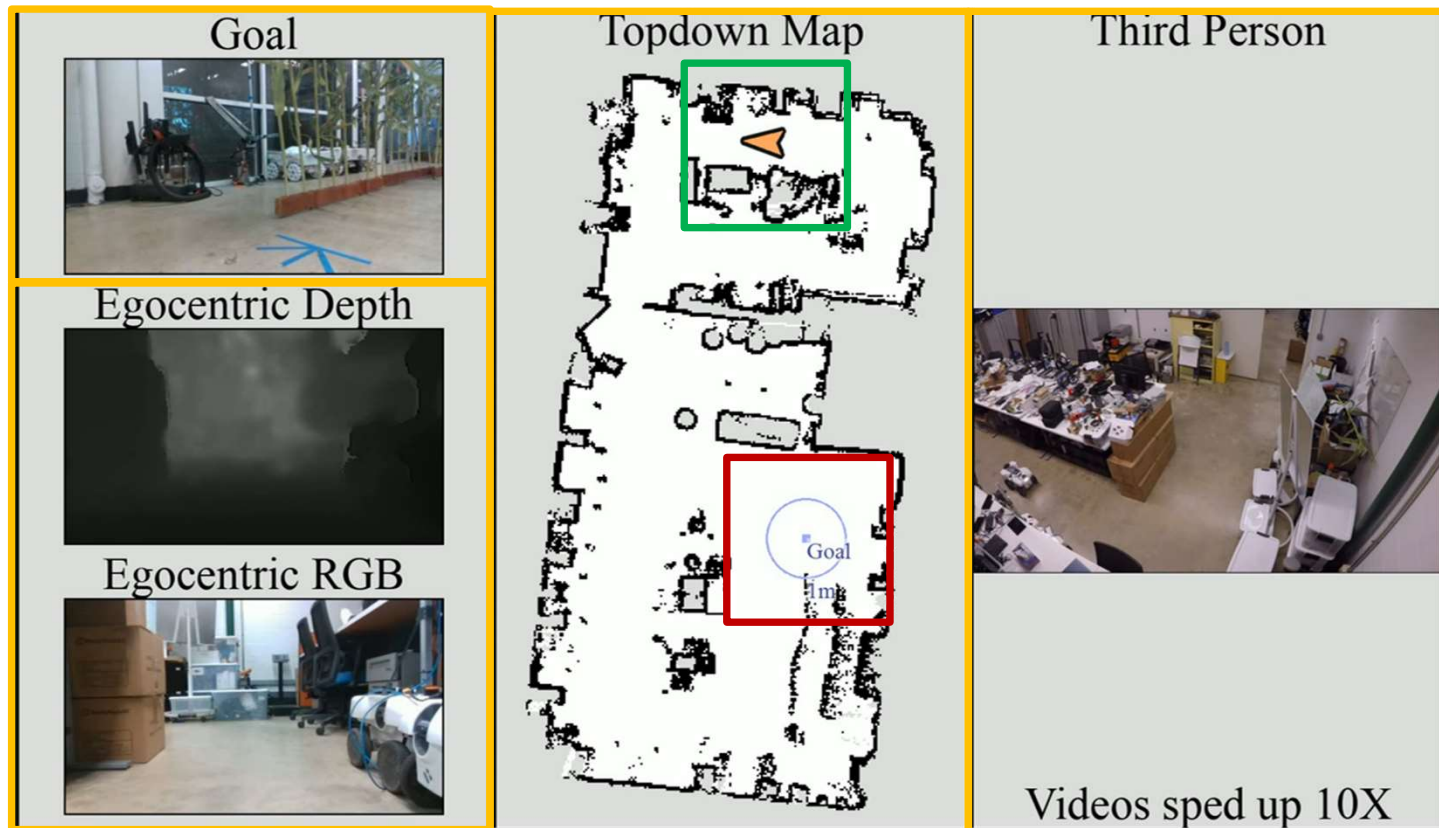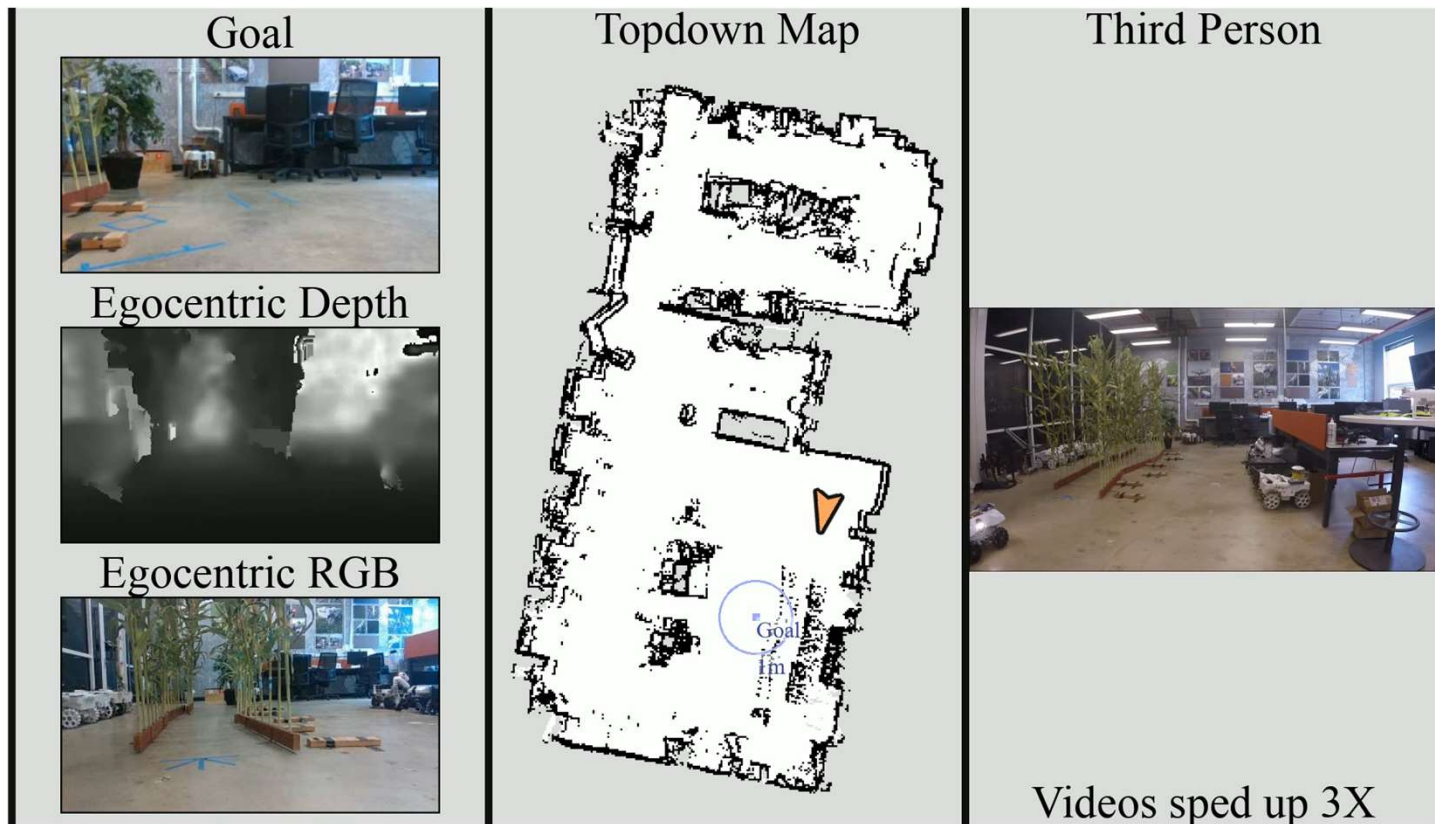(5 – 10m)

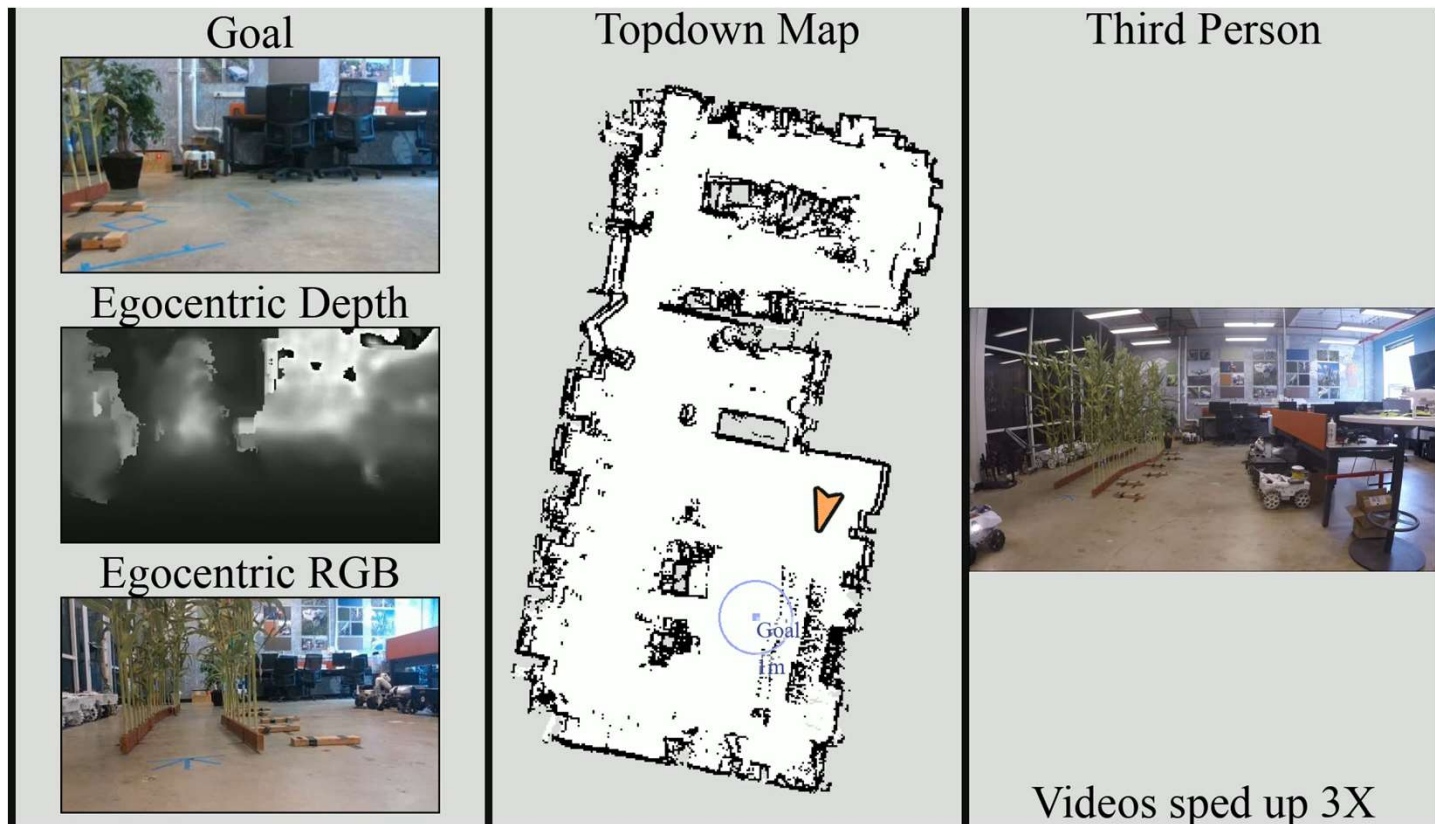# Image goal navigation in real life

# Image goal navigation in real life

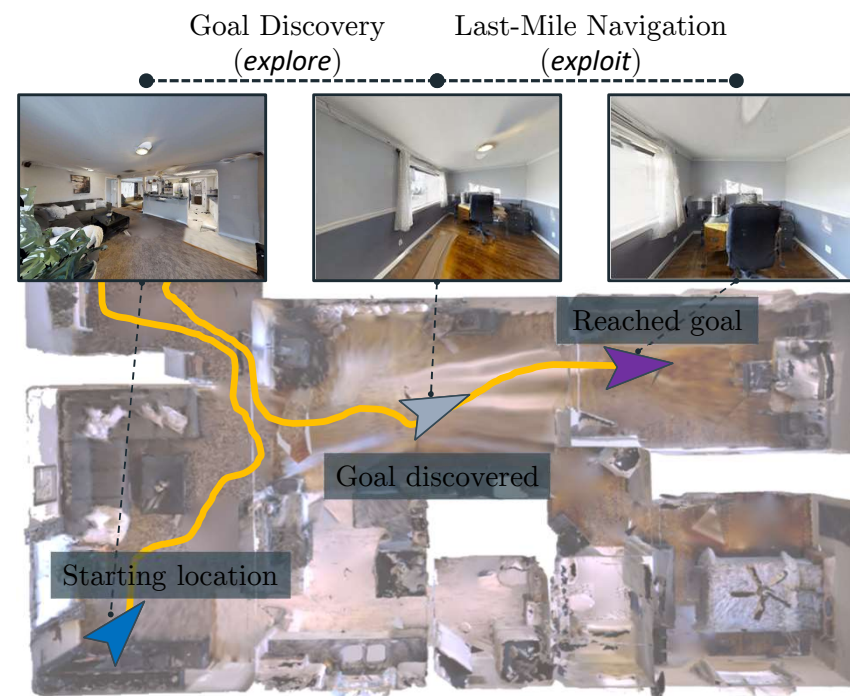# NRNS for the last mile

# SLING for the last mile
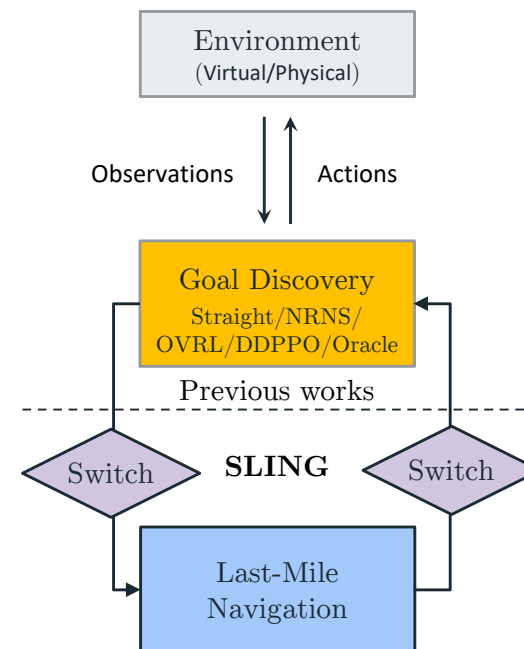
Last-Mile Embodied Visual Navigation

# Summary

# Summary

❖ **Goal discovery and last mile navigation**



Goal Discovery
(*explore*)

Last-Mile Navigation
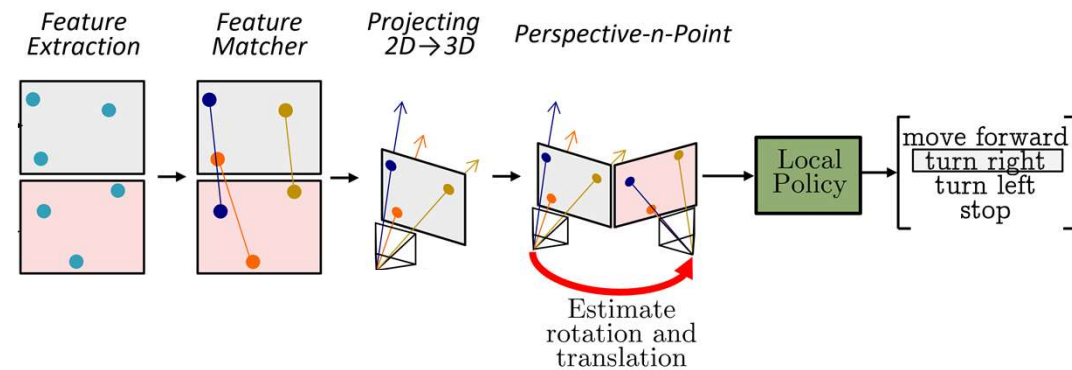(*exploit*)

Reached goal

Goal discovered

Starting location

# Summary

❖ Goal discovery and last mile navigation

❖ Switches to attach to prior baselines

# Summary

❖ Goal discovery and last mile navigation

❖ Switches to attach to prior baselines

❖ Principled 3D vision approach for last mile navigation

# Summary

❖ Goal discovery and last mile navigation

❖ Switches to attach to prior baselines

❖ Principled 3D vision approach for last mile navigation
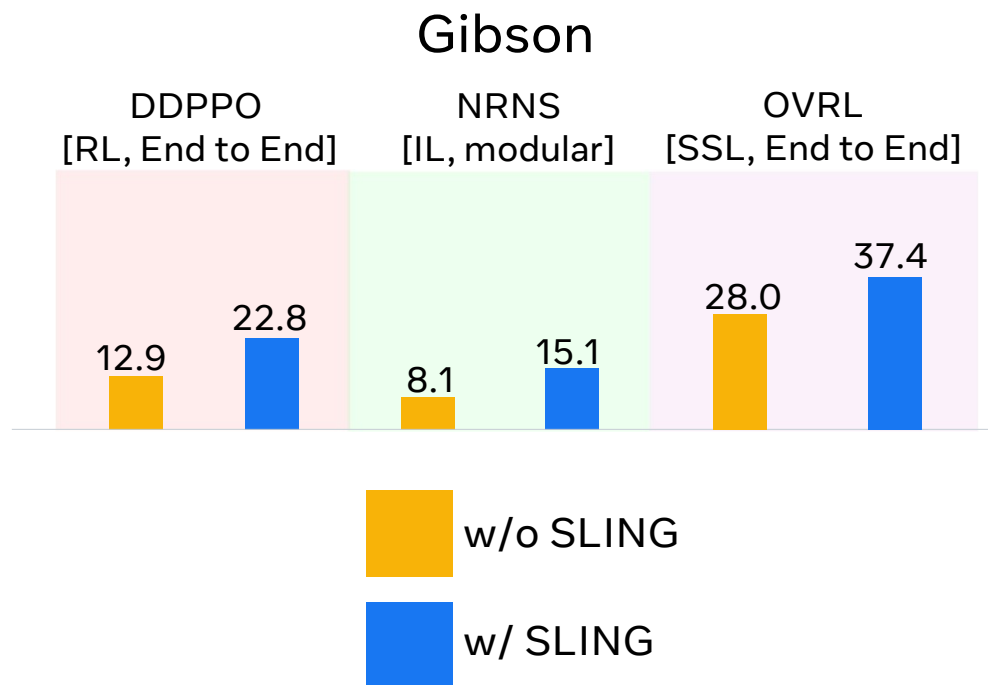
❖ Improved results across 5 baselines



Gibson

# Summary

❖ Goal discovery and last mile navigation

❖ Switches to attach to prior baselines

❖ Principled 3D vision approach for last mile navigation

❖ Improved results across 5 baselines

❖ Transfer to the real world