# Behavioral Analysis of Vision-and-Language Navigation Agents

Zijiao Yang
Oregon State University
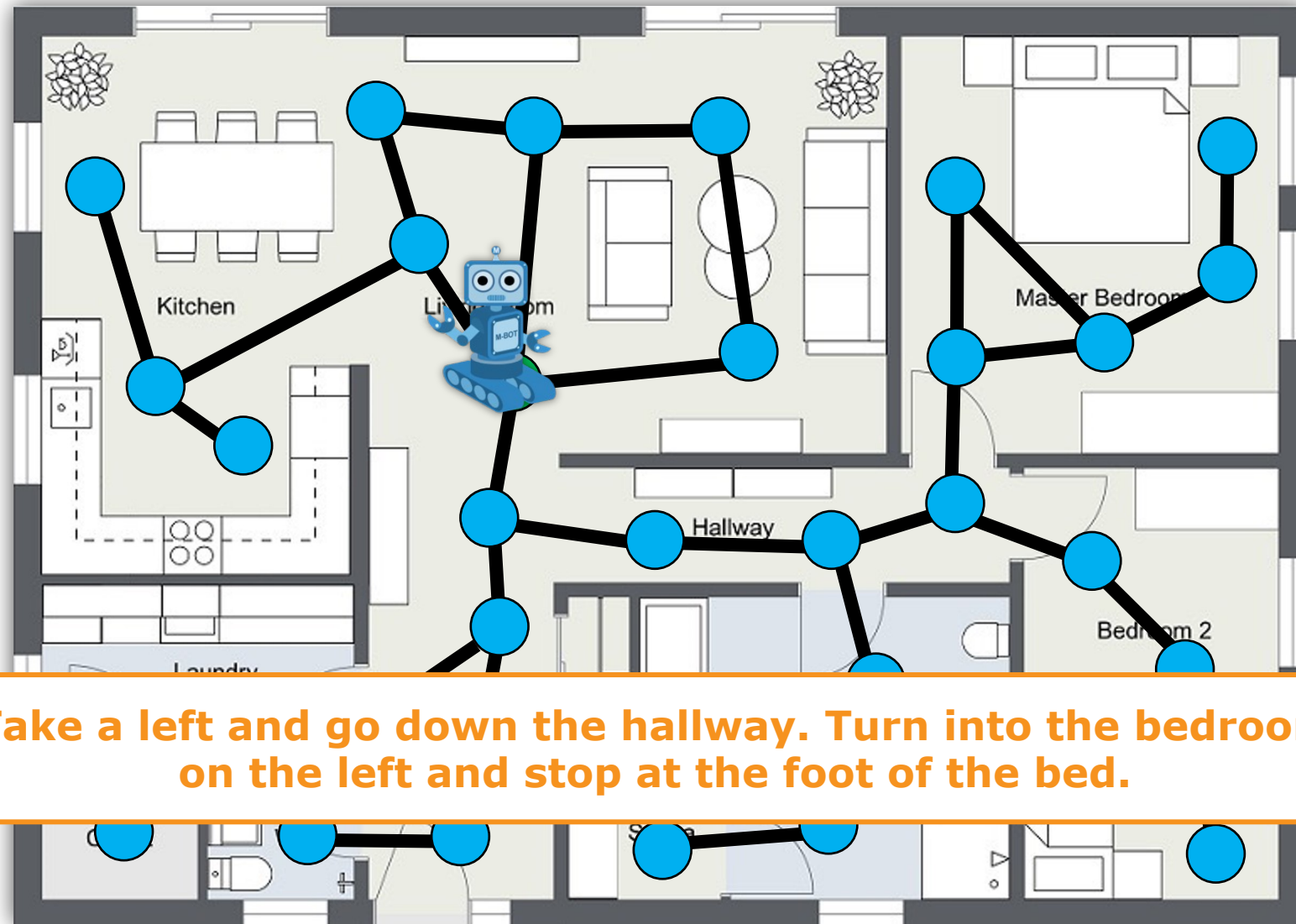
Arjun Majumdar
Georgia Institute of Technology

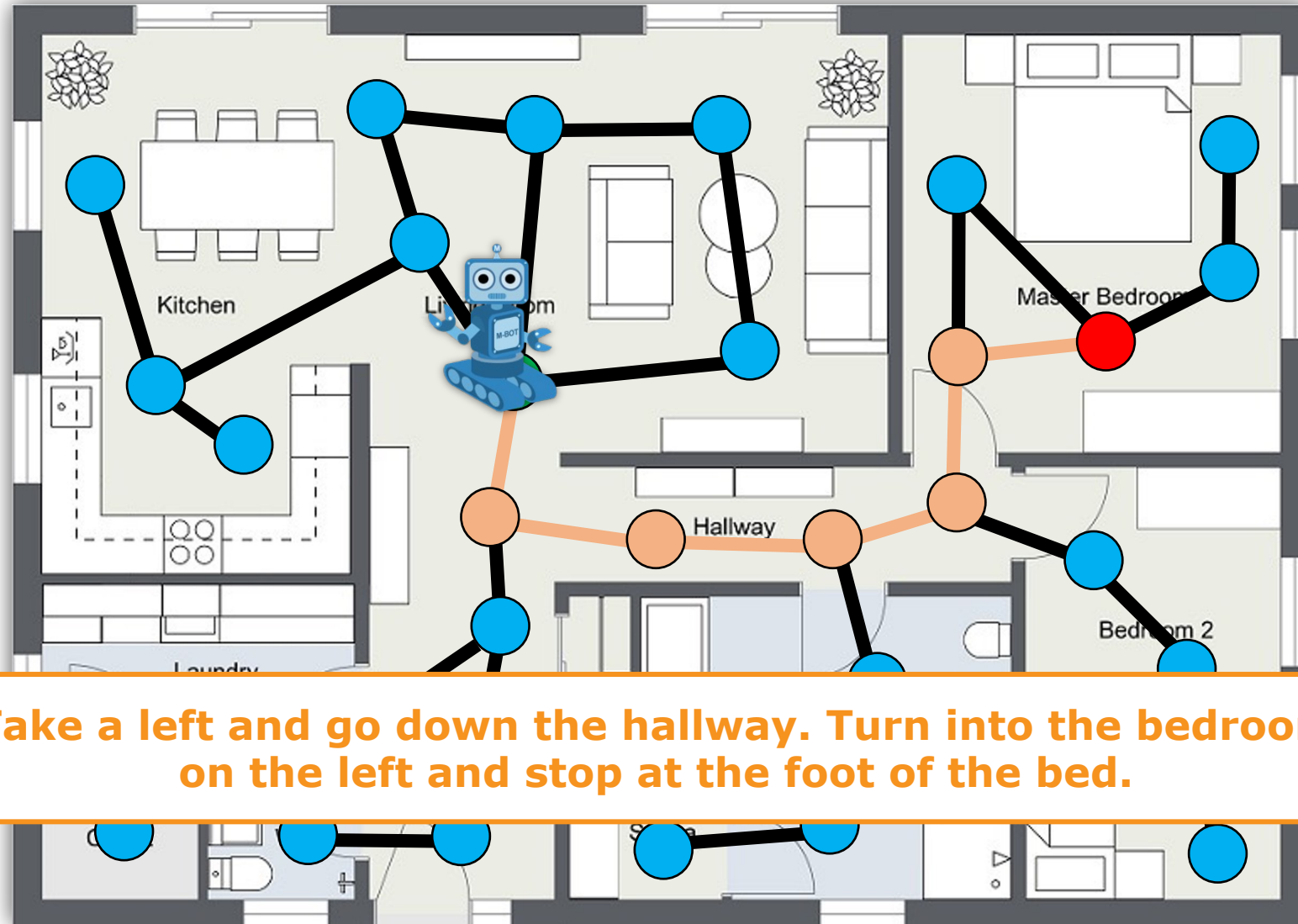Stefan Lee
Oregon State University
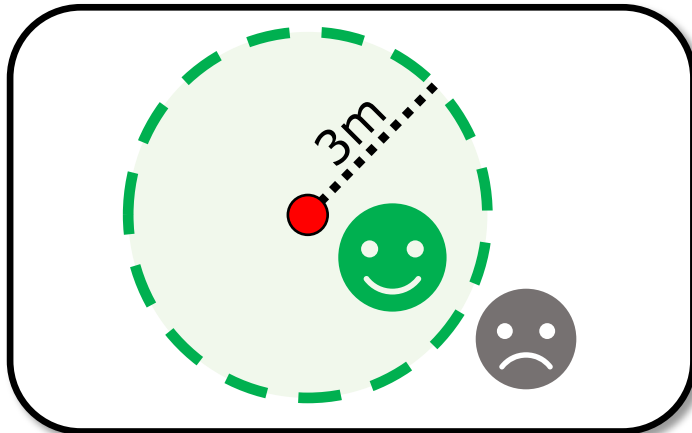
Session id: TUE-AM-245

Take a left and go down the hallway. Turn into the bedroom on the left and stop at the foot of the bed.
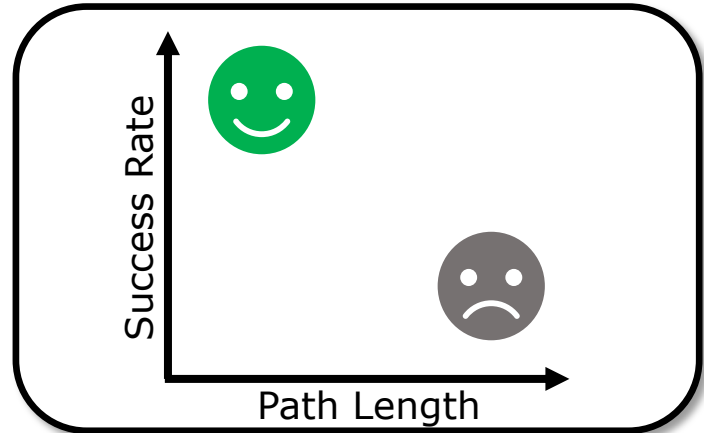
# Vision-and-Language Navigation (VLN)



Take a left and go down the hallway. Turn into the bedroom on the left and stop at the foot of the bed.
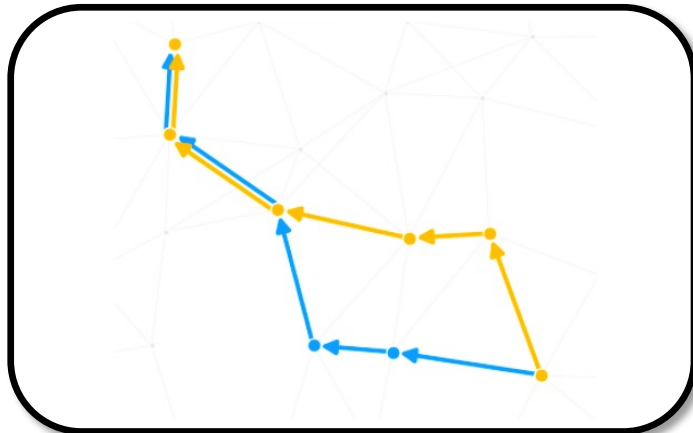
# Metrics measure aggregate performance



**Success Rate (SR)**

**Success weighted by Path Length (SPL)**

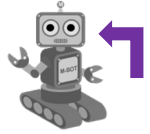**Normalized Dynamic Time Warping (nDTW)**

Useful for comparing methods in aggregate but don't speak to fine-grained agent capabilities.

# Examine Four Fine-grained Skill Competencies

Explicit Stop                                    *"This is your destination"*

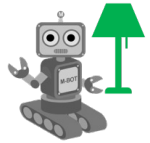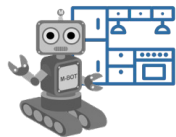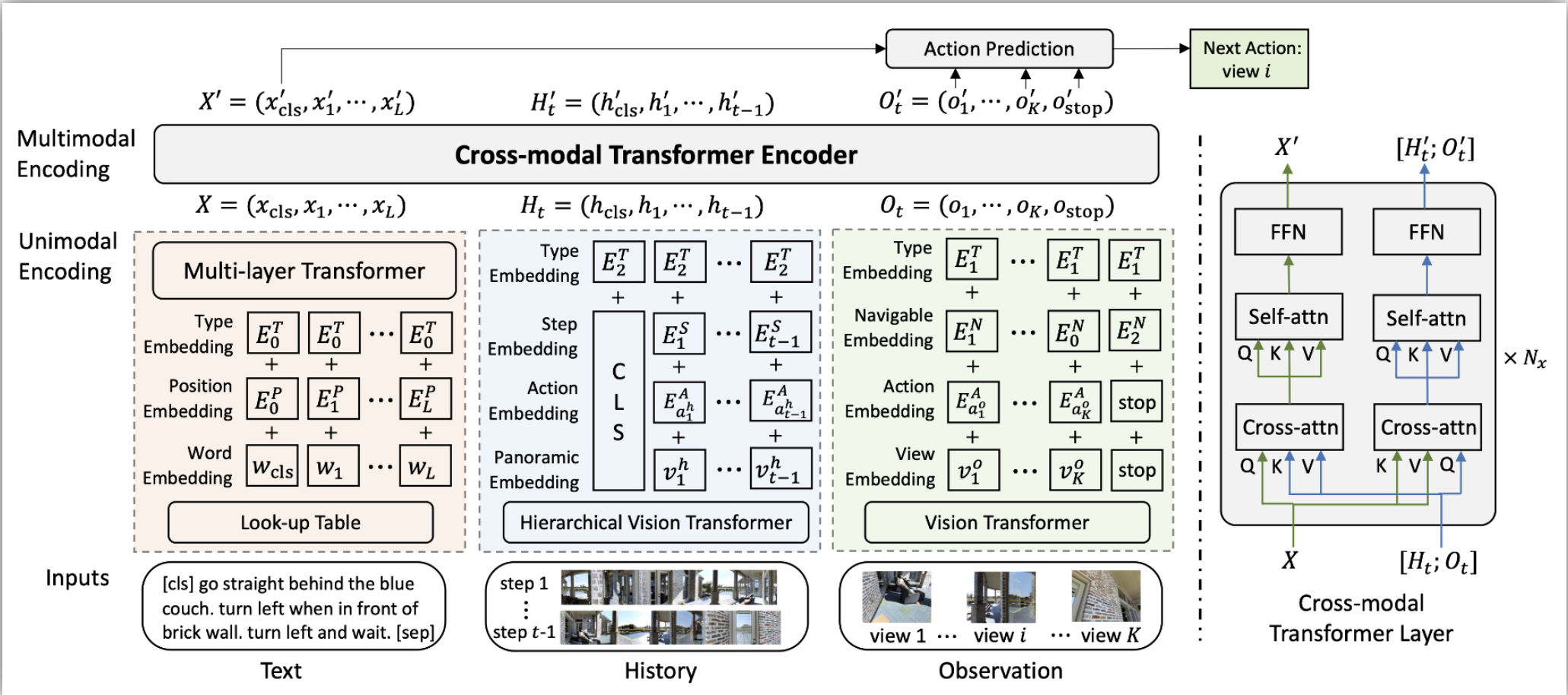Unconditional Directional                        *"Make a left"*

Object Seeking                                   *"Go towards the lamp"*

Room Seeking                                     *"Go to the kitchen"*
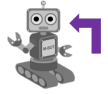
# Case Study



*History Aware Multimodal Transformer for Vision-and-Language Navigation (Chen et al.)*

# Summary of Findings

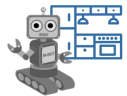Explicit Stop       *"This is your destination"*       ✔

Unconditional Directional       *"Make a left"*       ✔
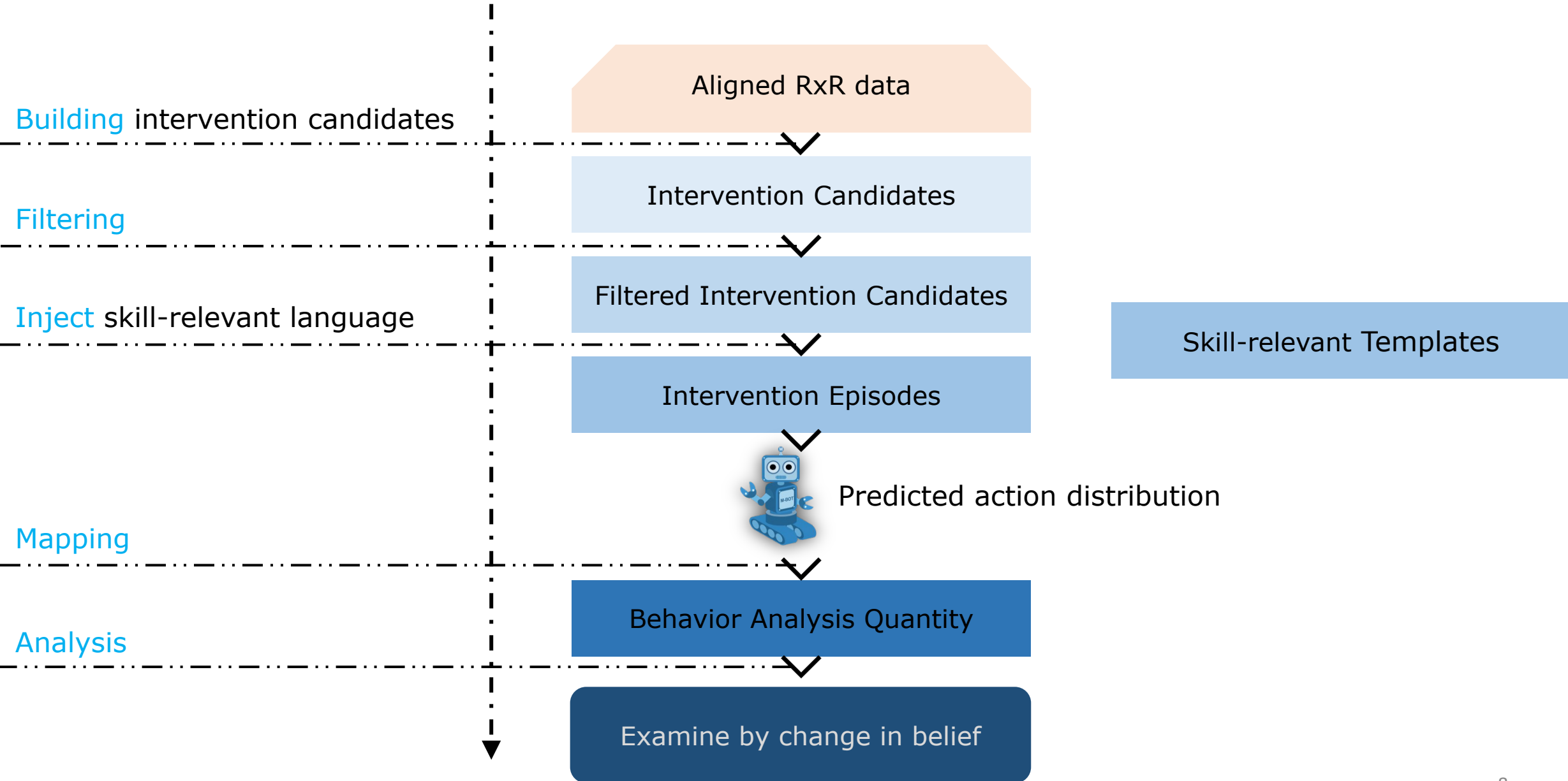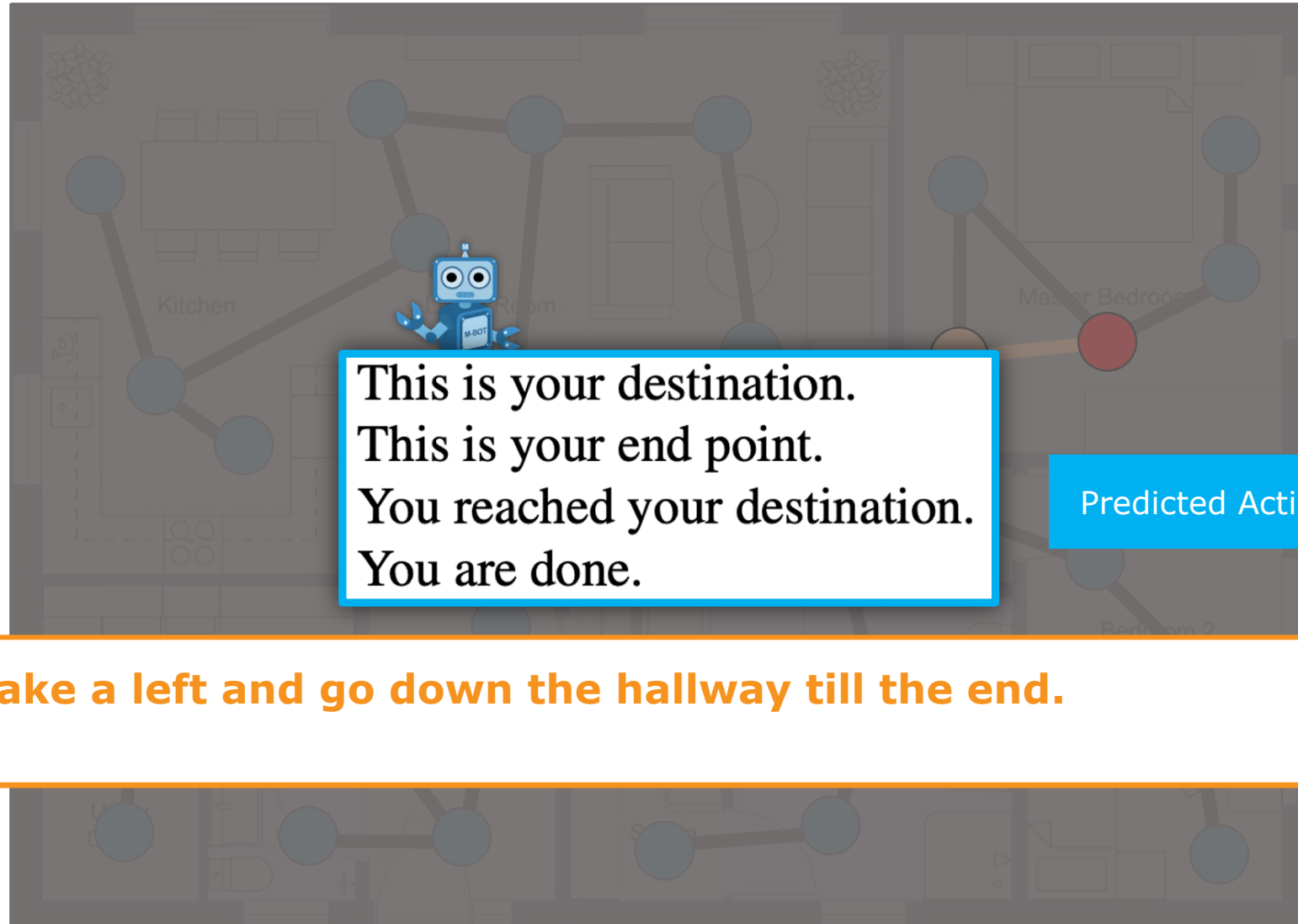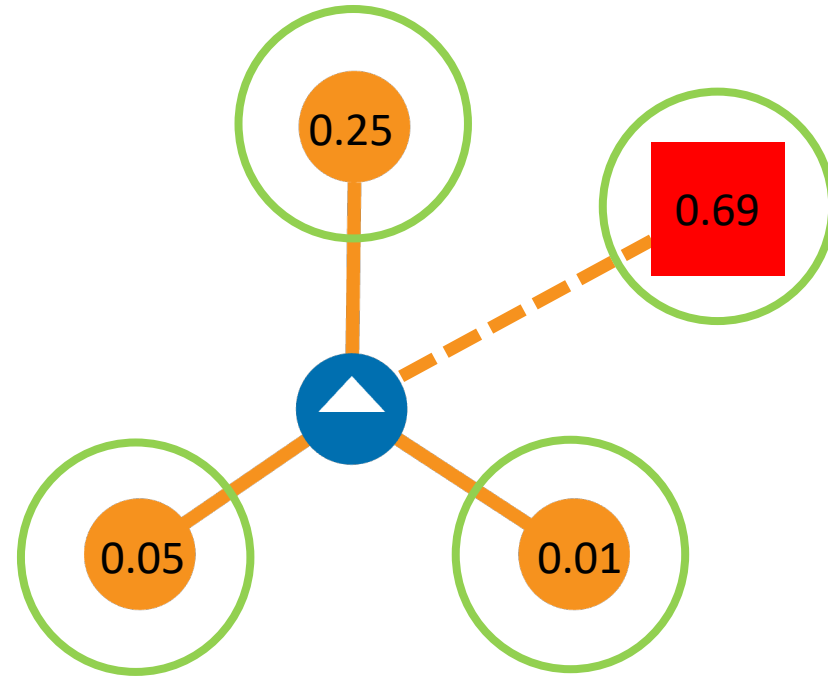
Object Seeking       *"Go towards the lamp"*       ?

Room Seeking       *"Go to the kitchen"*       ?

# Behavioral Analysis Framework



Building intervention candidates

Filtering

Inject skill-relevant language

Mapping

Analysis

Aligned RxR data

Intervention Candidates

Filtered Intervention Candidates

Skill-relevant Templates

Intervention Episodes

Predicted action distribution

Behavior Analysis Quantity

Examine by change in belief

# HAMT: Stop (Explicit Stop)



This is your destination.
This is your end point.
You reached your destination.
You are done.
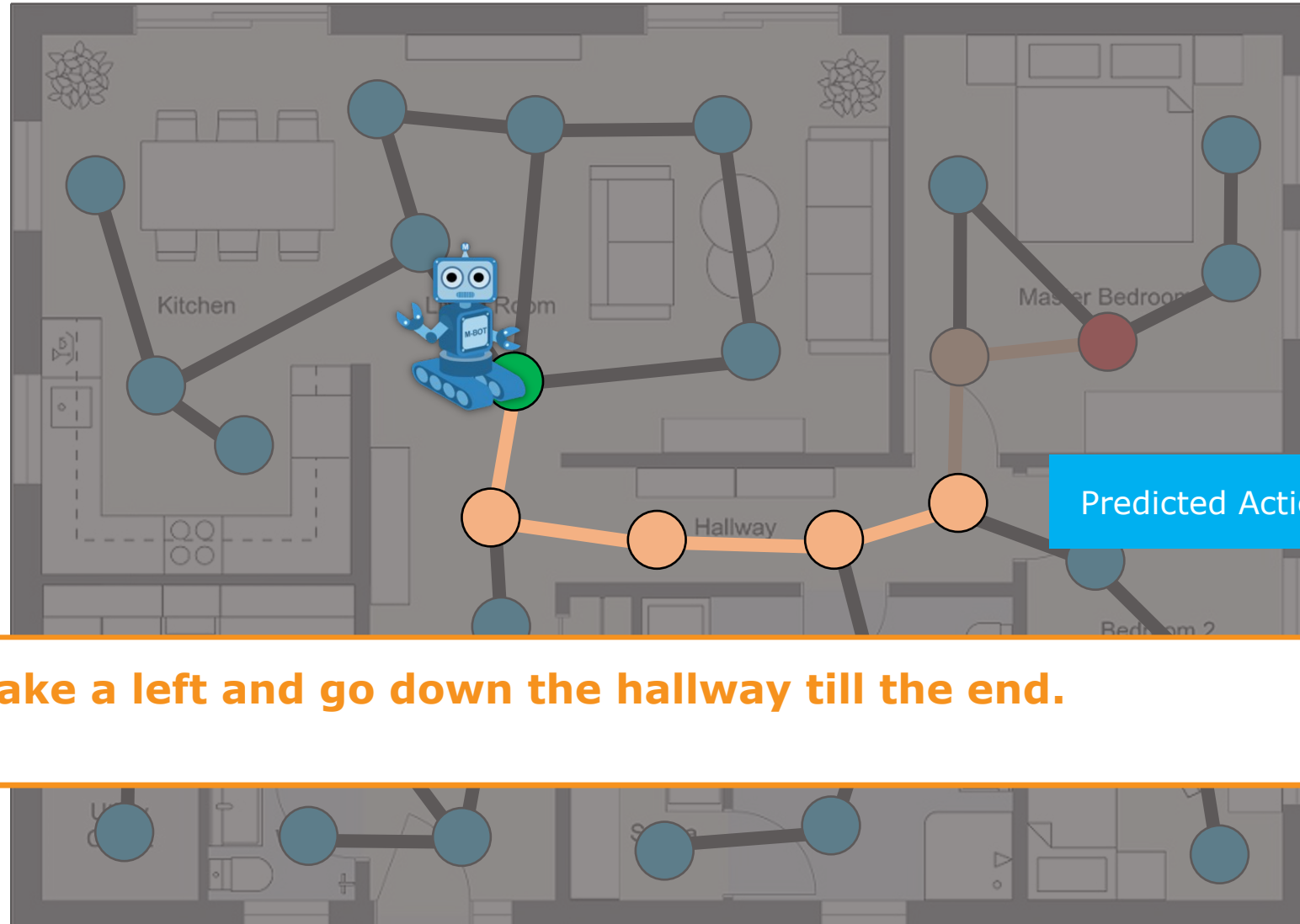
Predicted Action Distribution

**Take a left and go down the hallway till the end.**
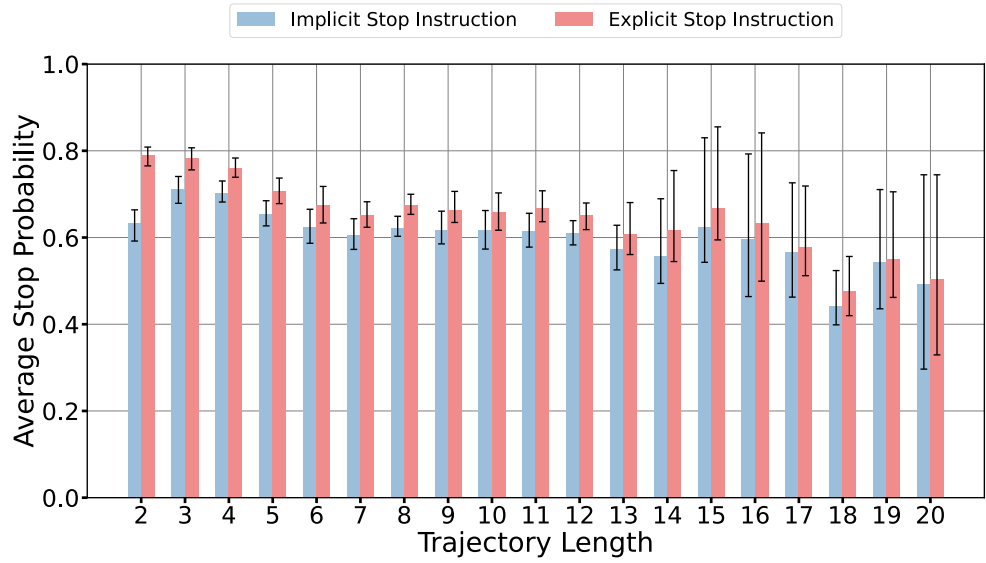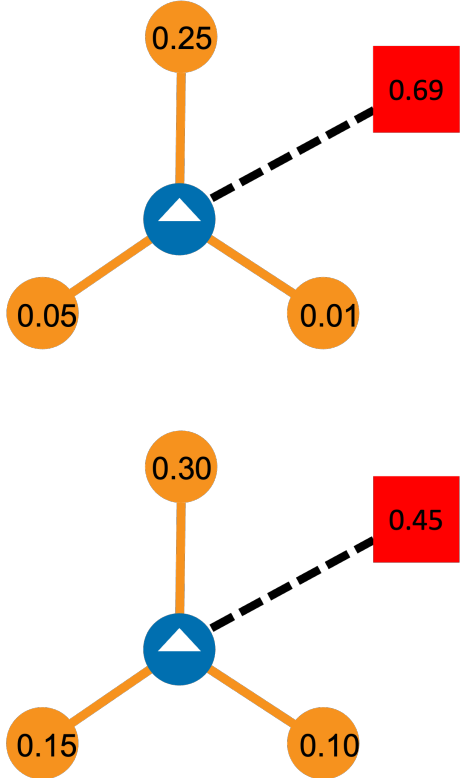
# HAMT: Explicit Stop

# HAMT: Stop (Implicit Stop)



Predicted Action Distribution

**Take a left and go down the hallway till the end.**

# HAMT: Stop (Analysis)

# HAMT: Stop (Analysis)

# HAMT: Unconditional Directional Instructions (Filtering)

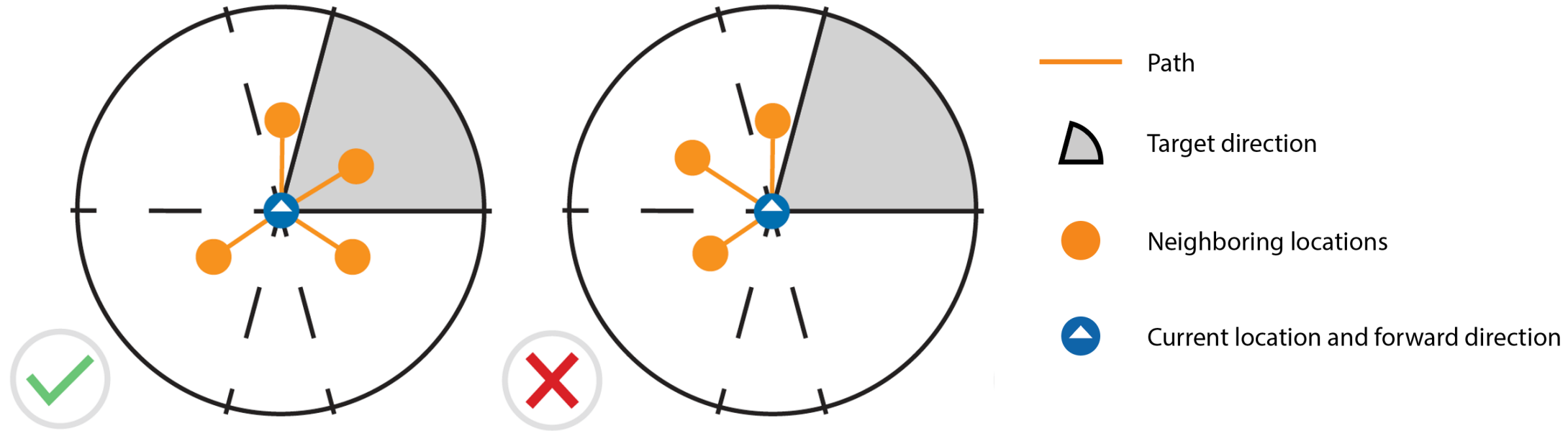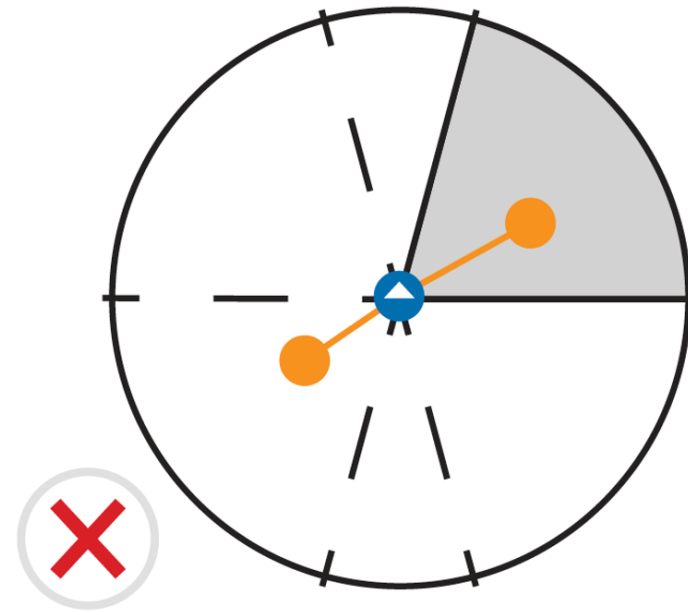# HAMT: Unconditional Directional Instructions (Filtering)



Path

Target direction

Neighboring locations

Current location and forward direction

# HAMT: Unconditional Directional Instructions (Filtering)



Path

Target direction

Neighboring locations

Current location and forward direction

# HAMT: Unconditional Directional Instructions (Templates)

Template:

Walk forward. (*forward*)
Turn around and walk forward. (*backward*)
Turn left and walk forward. (*left*)
Turn right and walk forward. (*right*)
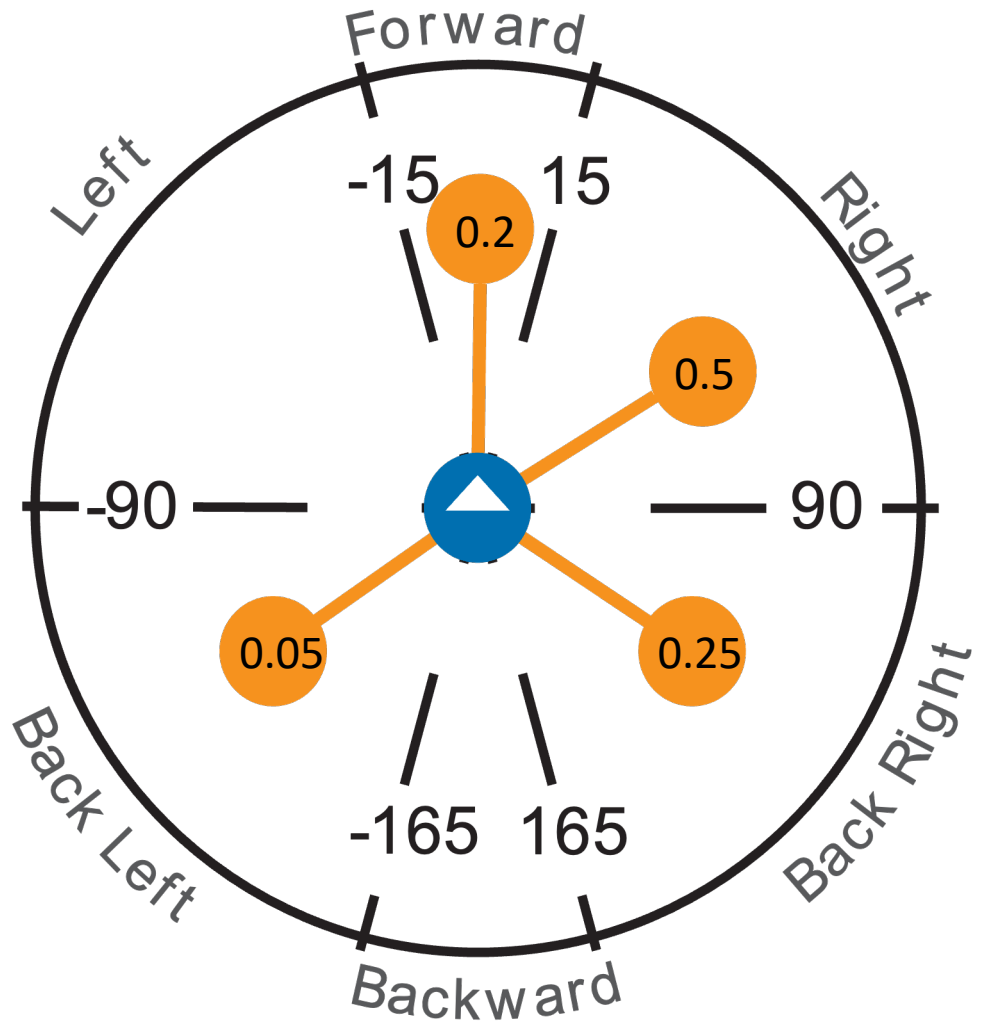Turn around and go to your right. (*back left*)
Turn around and go to your left. (*back right*)

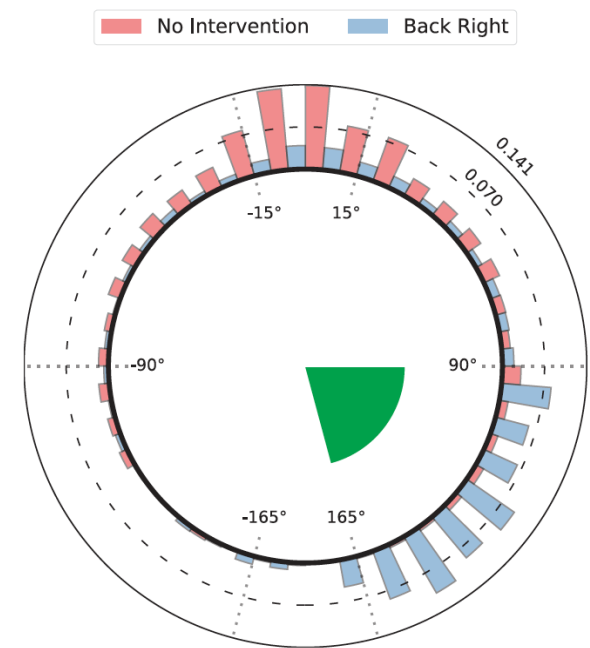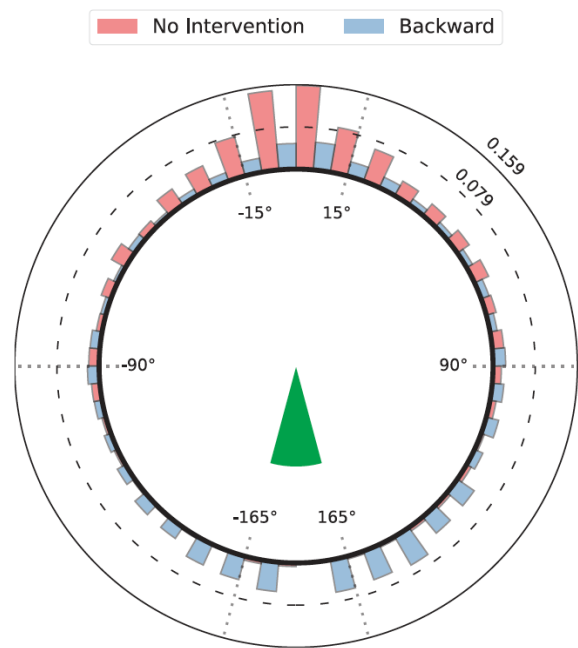# HAMT: Unconditional Directional Instructions (Mapping)

# HAMT: Unconditional Directional Instructions (Analysis)

# HAMT: Object Seeking (Filtering)

REVERIE (Qi et al.)

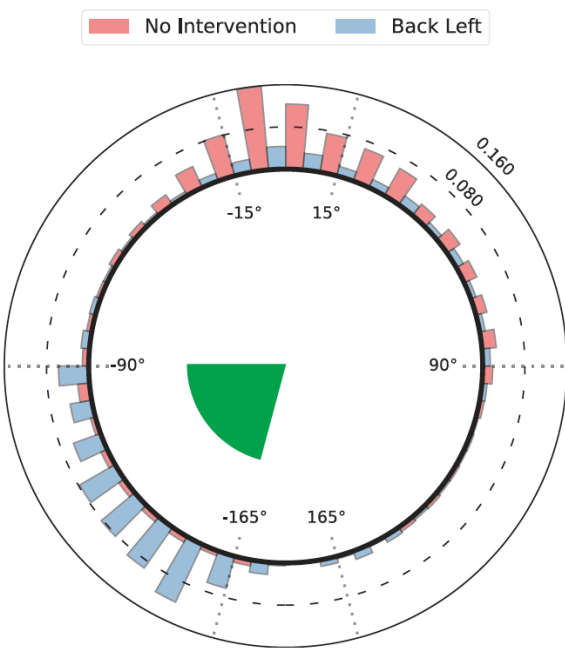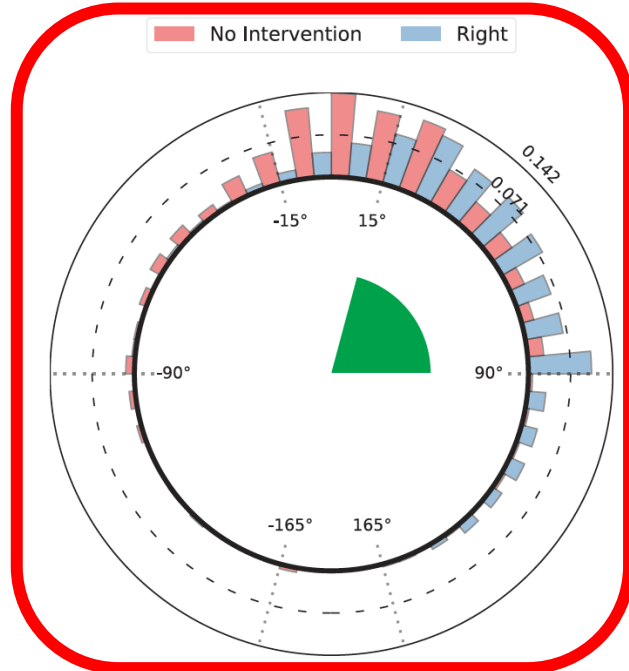Walk towards the XX (*Object*)



15°
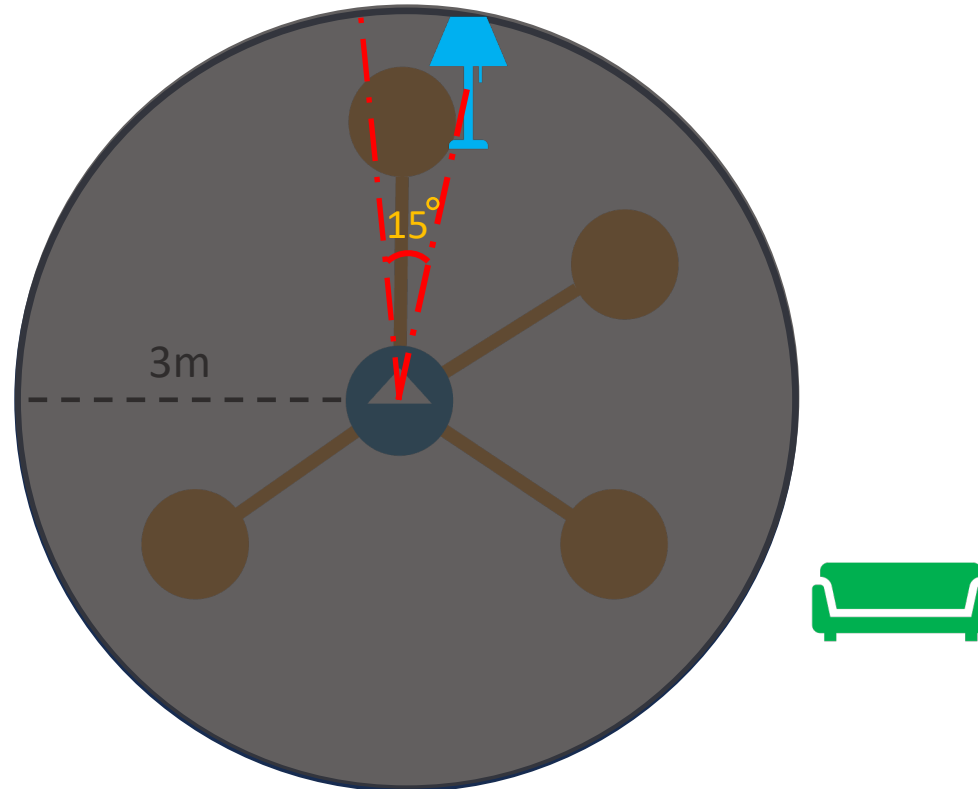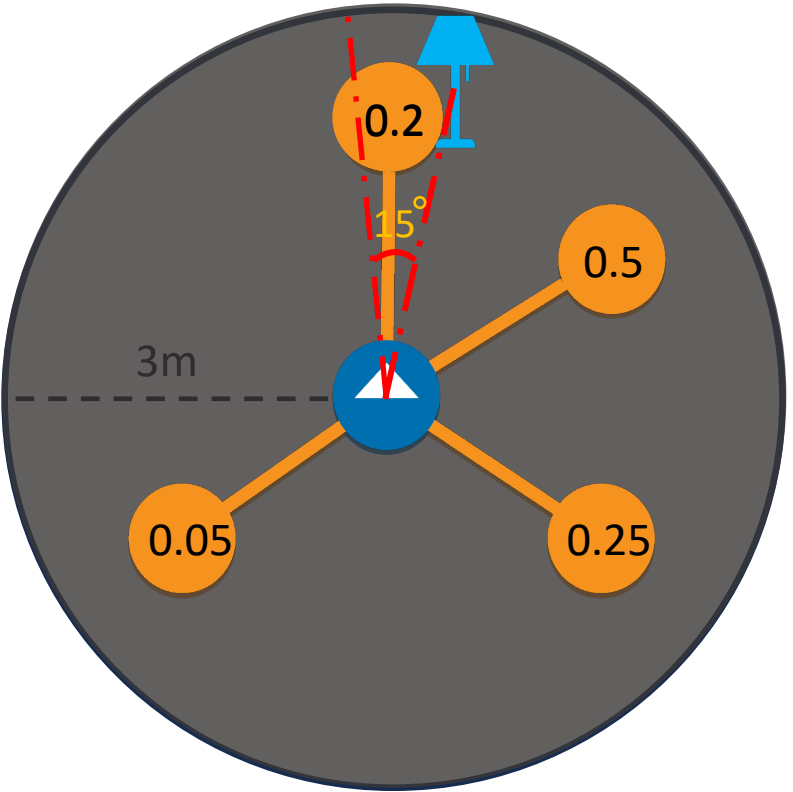
3m

# HAMT: Object Seeking (Mapping)
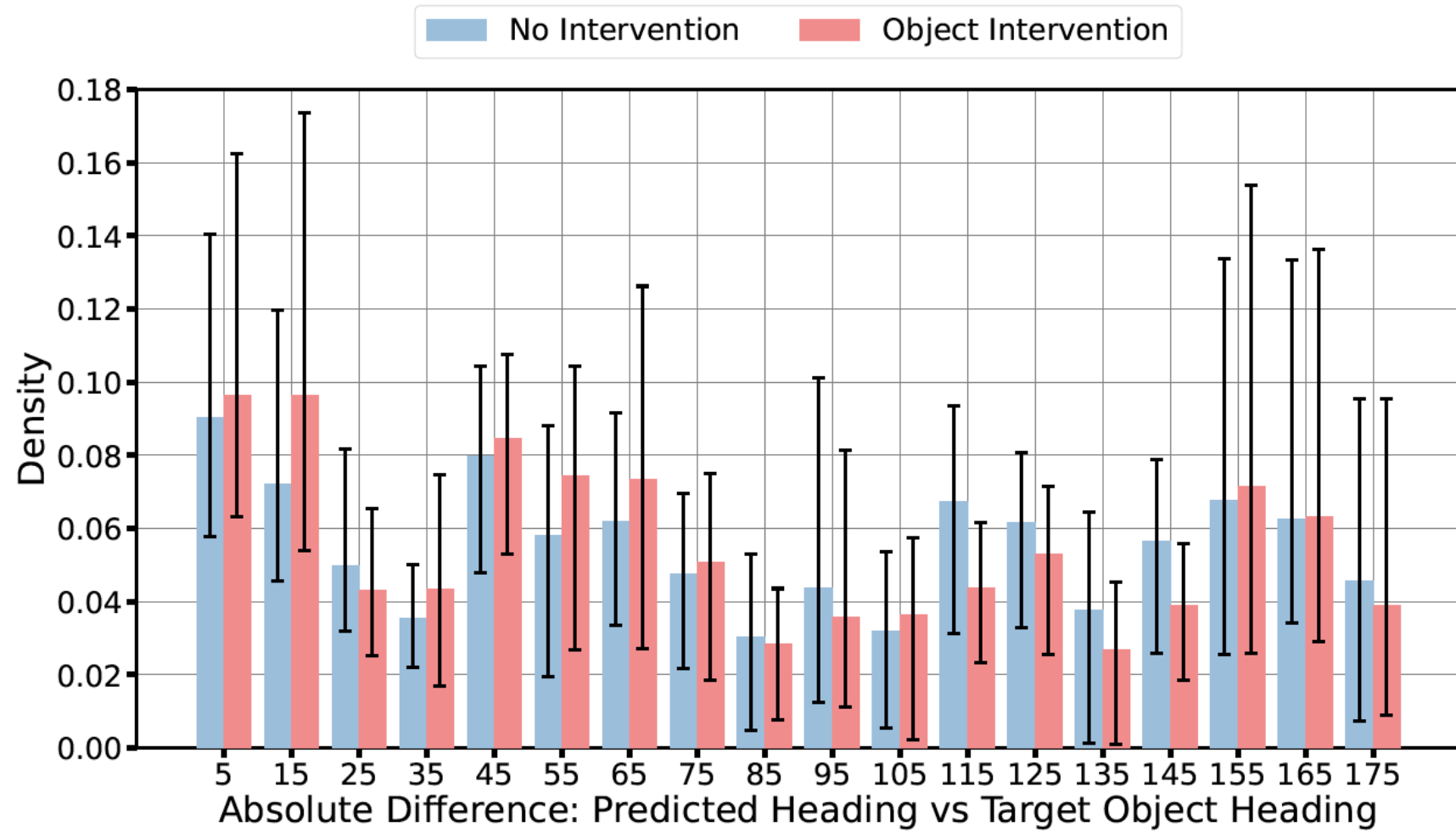
REVERIE (Qi et al.)
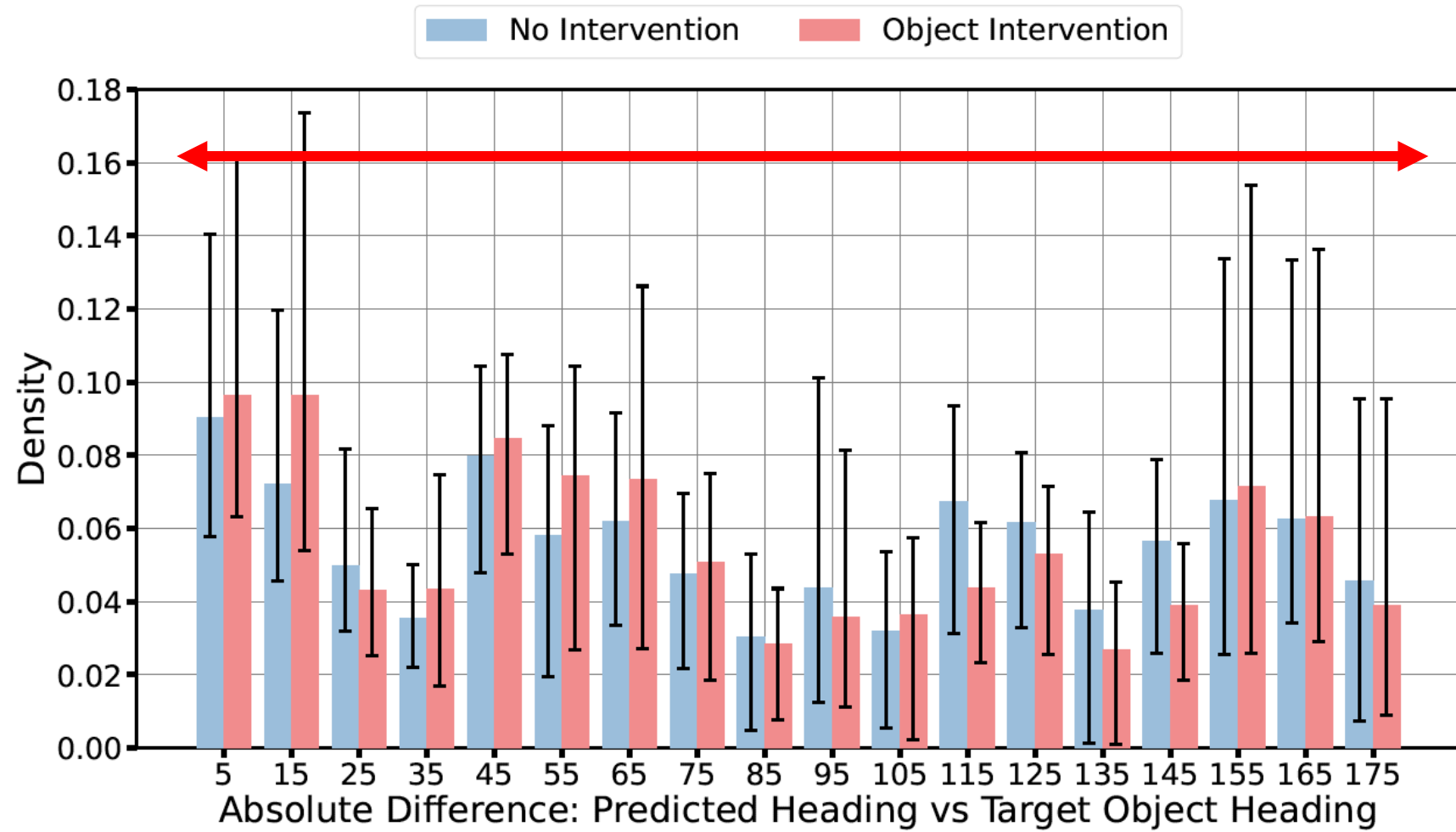


Behavior Analysis Quantity ➡ Absolute angular error relative to object
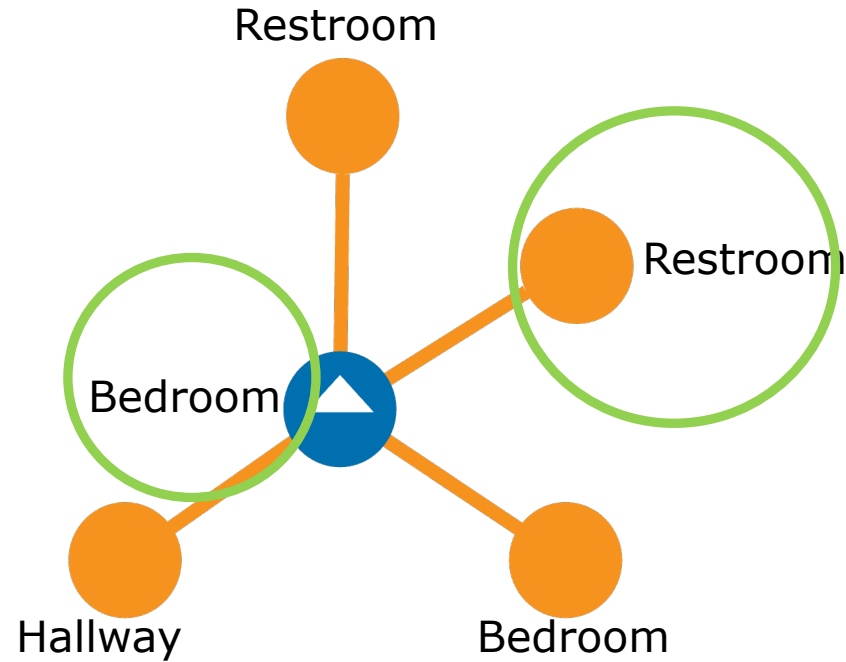
# HAMT: Object Seeking (Analysis)

# HAMT: Object Seeking (Analysis)

# HAMT: Room Seeking (Filtering)

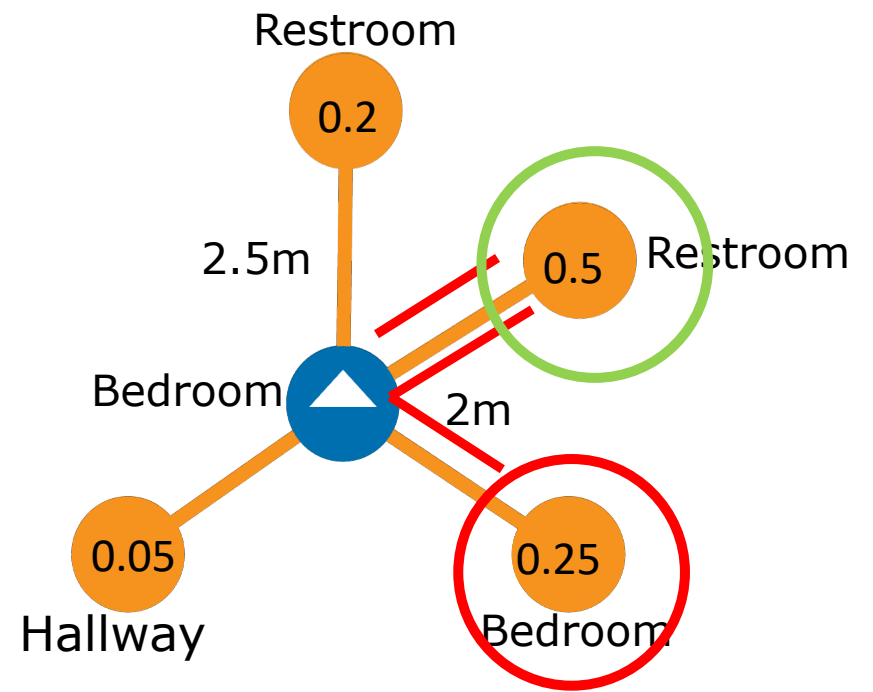MATTERPORT3D (Chang et al.)



Walk towards the XX (*room type*)

# HAMT: Room Seeking (Mapping)
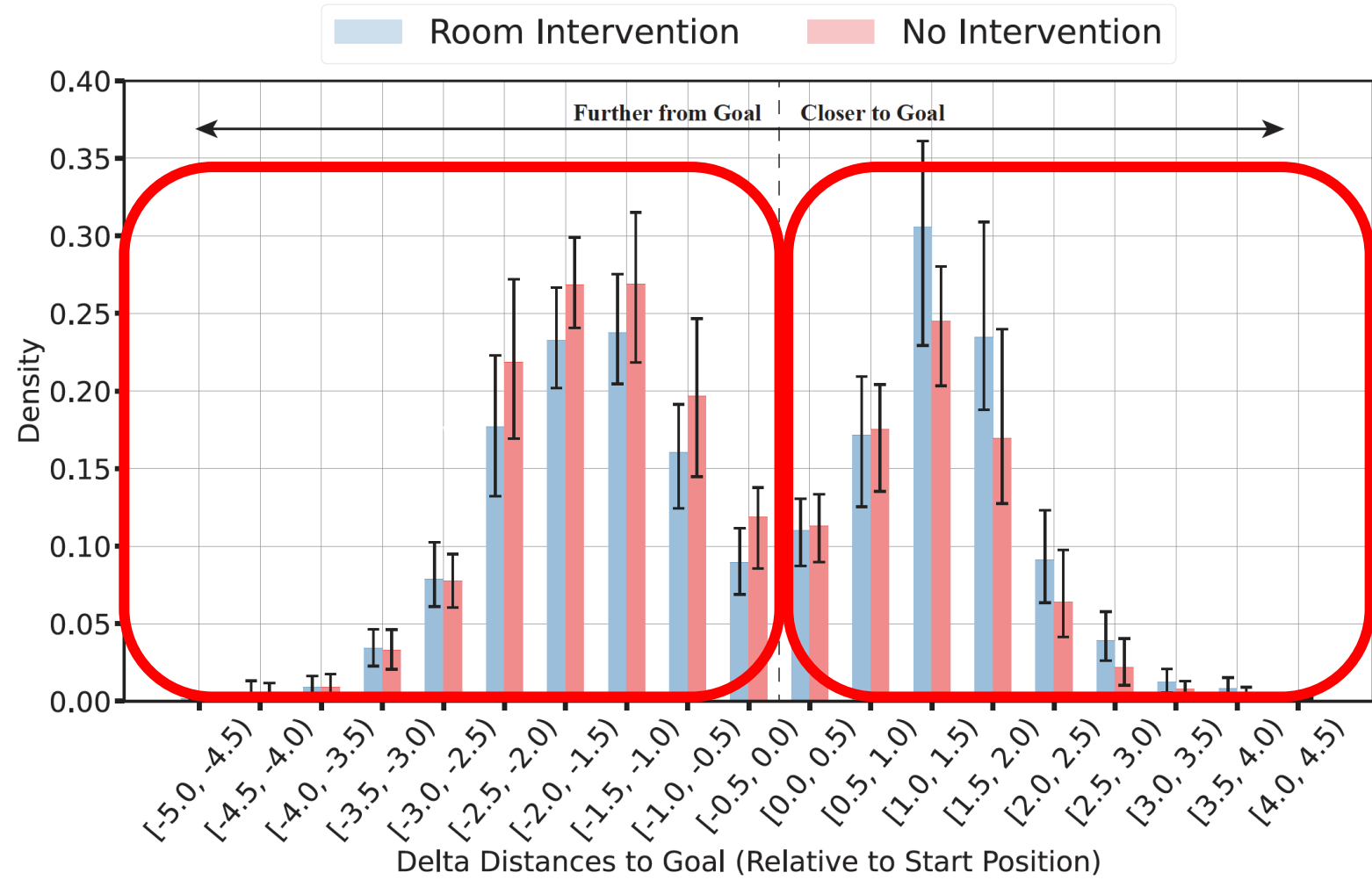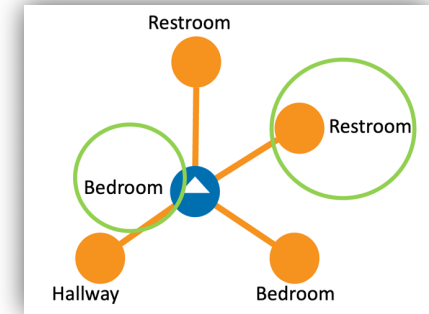
MATTERPORT3D (Chang et al.)

Length( ) ── Length( ) :



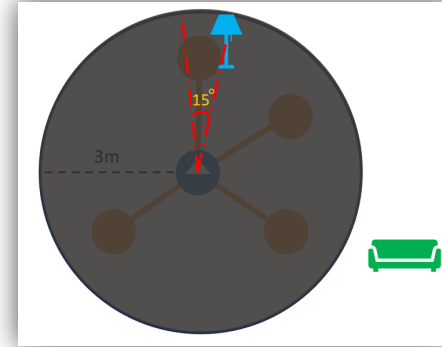Walk towards the restroom.

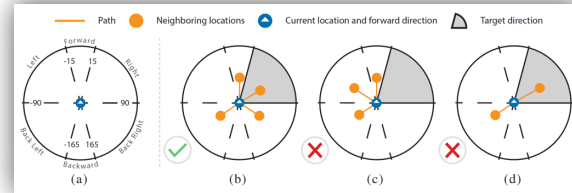# HAMT: Room Seeking (Analysis)

# Competency Score



Target specified by Intervention Instruction

# Competency Score

Ordered by increasing RxR performance

| Method | Stop | Turn | Object | Room | Avg. |
|---|---|---|---|---|---|
| EnvDrop | 62.65 | 27.14 | 11.06 | 23.64 | 31.12 |
| EnvDrop (ViL CLIP) | 66.76 | 27.45 | 12.83 | 26.82 | 33.47 |
| HAMT | 71.65 | 43.74 | 12.00 | 26.63 | 38.50 |

# Competency Score

Ordered by increasing RxR performance

| Method | Stop | Turn | Object | Room | Avg. |
|---|---|---|---|---|---|
| EnvDrop | 62.65 | 27.14 | 11.06 | 23.64 | 31.12 |
| EnvDrop (ViL CLIP) | 66.76 | 27.45 | 12.83 | 26.82 | 33.47 |
| HAMT | 71.65 | 43.74 | 12.00 | 26.63 | 38.50 |

# Contact & Project Information

### GitHub Repository

https://github.com/Yoark/vln-behave

### Project Page

https://yoark.github.io/paper_sites/vln-behave/

*For more information or any questions, feel free to reach out:*
Zijiao Yang
yangziji@oregonstate.edu