

Representations for Visual Navigation and How to Train Them

Saurabh Gupta
UIUC

In this talk,

Representations for Places that Afford Navigation in Novel Environments

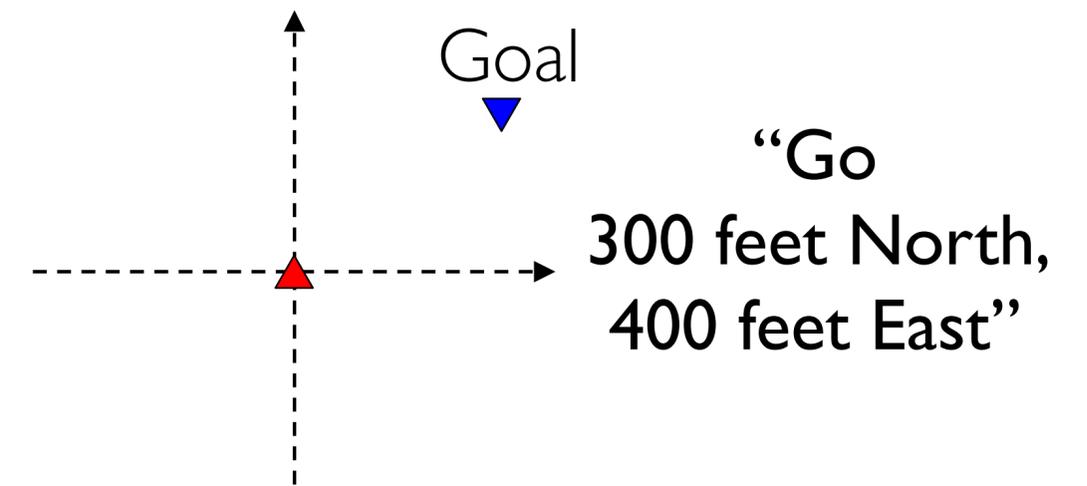
Basic Navigation Problems



Robot with a first person camera



Dropped into a novel environment

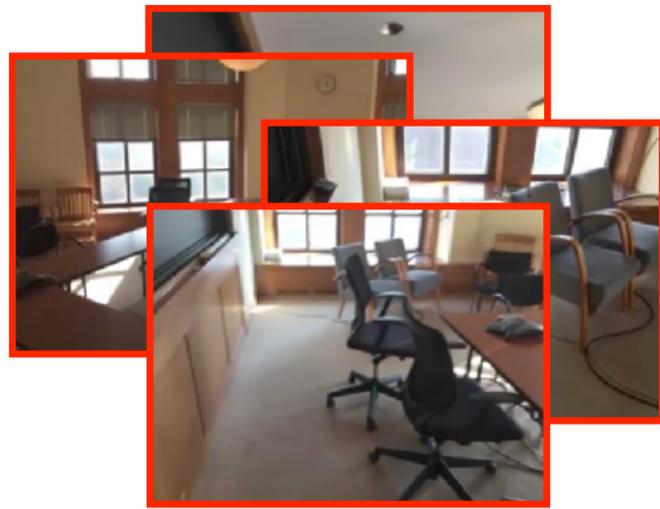


“Go Find a Chair”

“Explore the Environment”

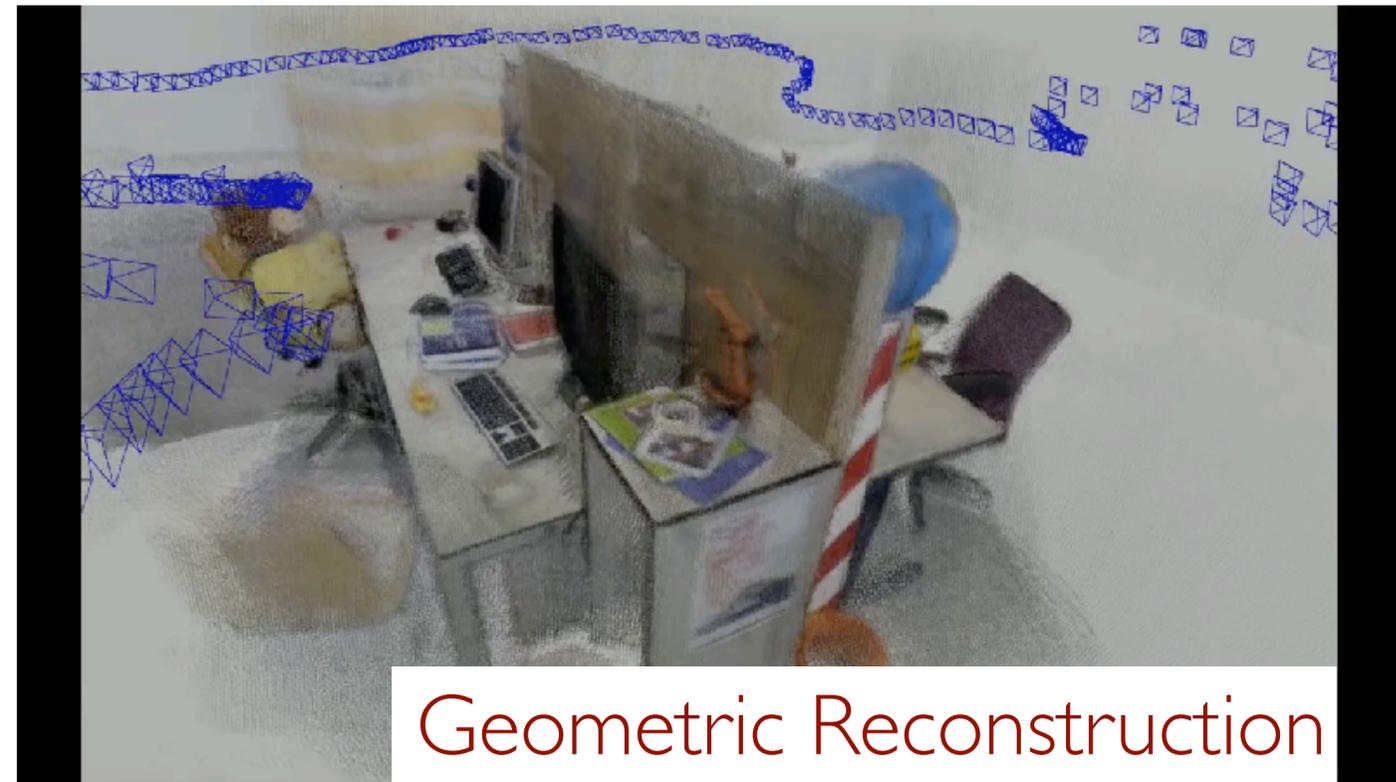
Discover paths or Explore

Classical Solution

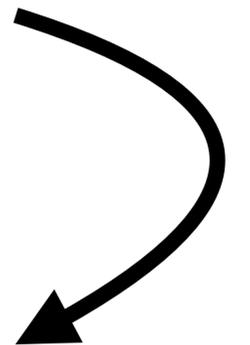


Observed Images

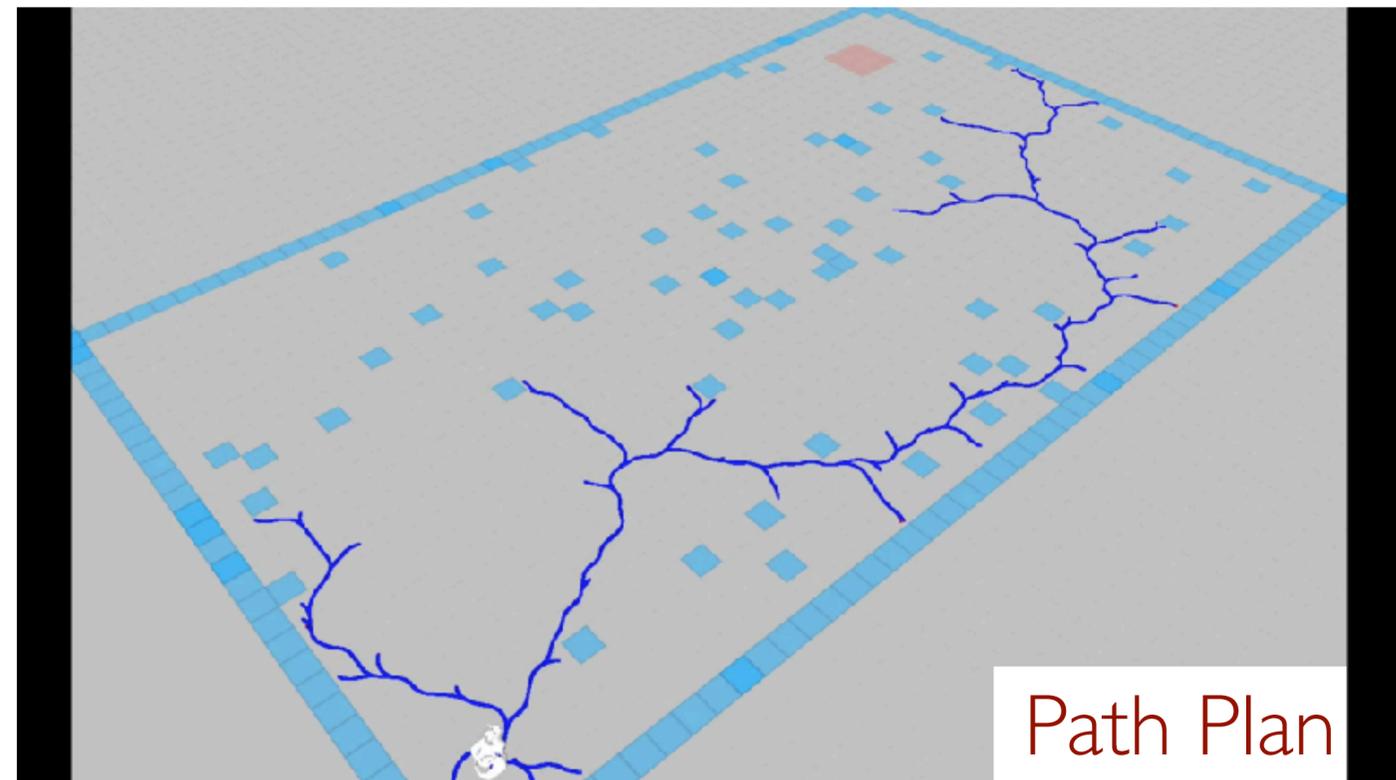
Mapping



Geometric Reconstruction



Planning



Path Plan

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. RAI 1996. 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 3D Reconstruction of the World

Unnecessary

Do we need to tediously reconstruct everything on this table?



Humans can do quite a bit without accurate metric 3D information

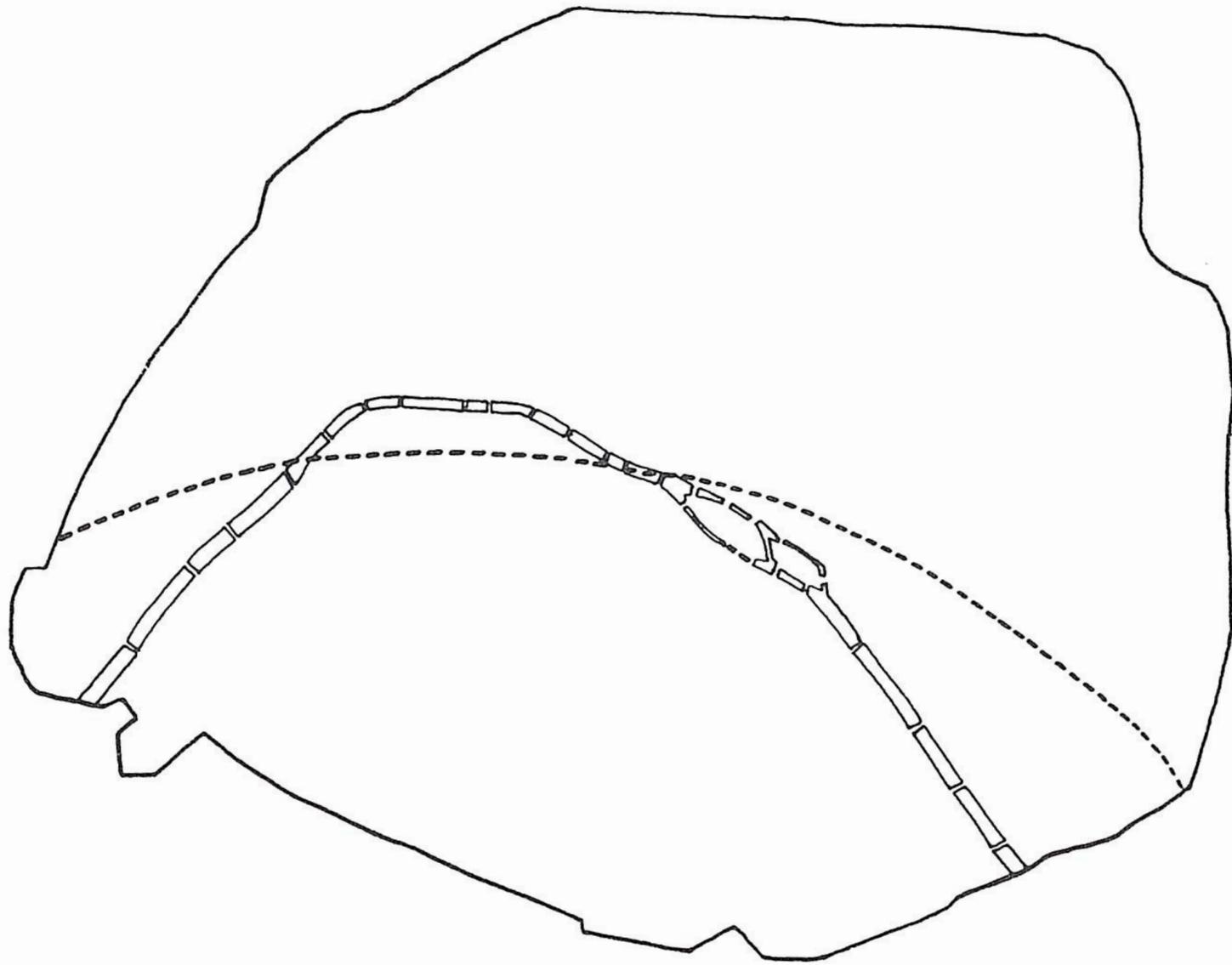


FIGURE 8.5
Perceived curvature of the Seine. The dotted line represents the median curvature imparted to the Seine in the subject's handdrawn maps. It is superimposed on the actual course of the river.

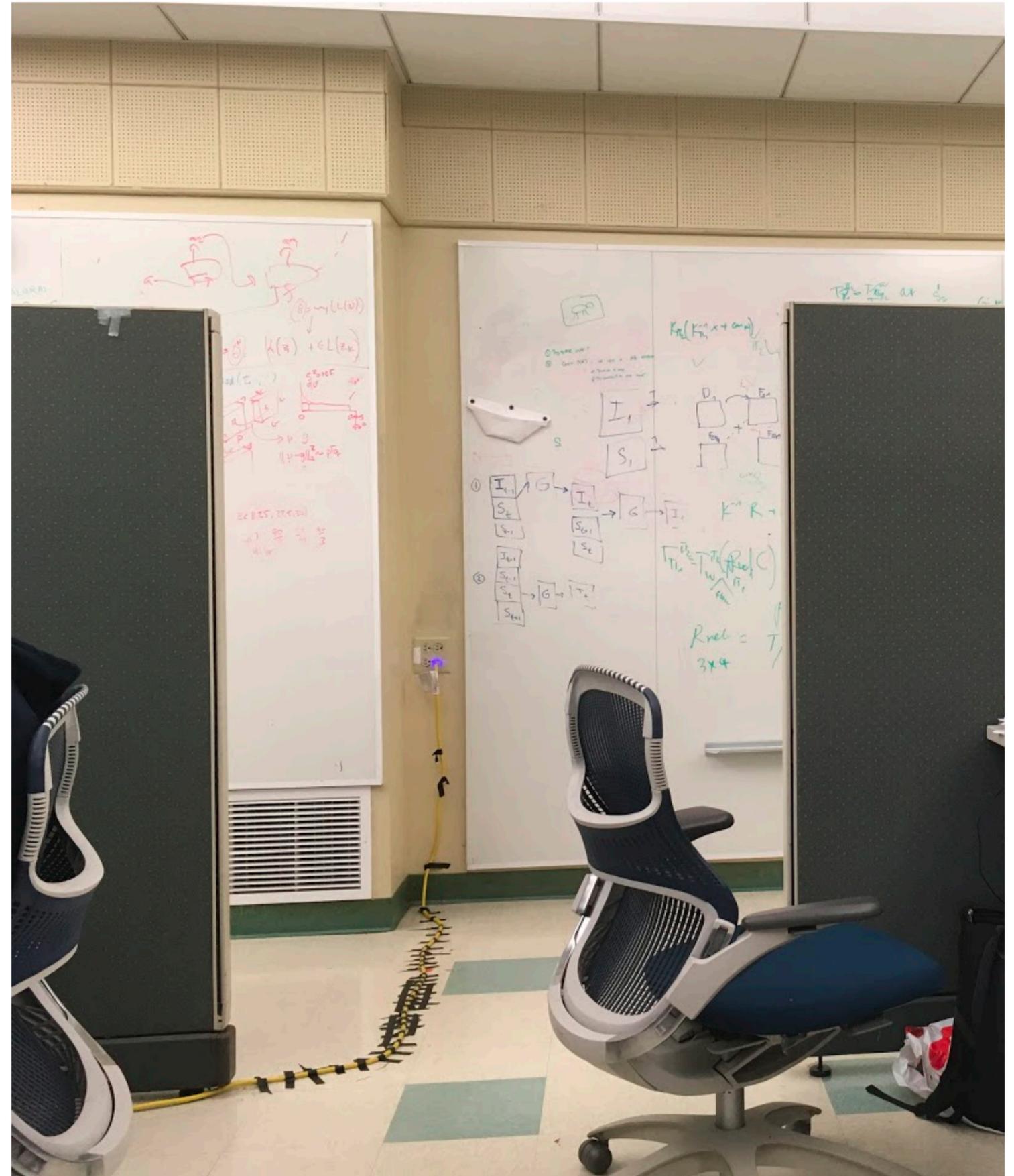
and much of the subjects subsequently...
But there is a serious distortion in the way the Seine is represented. In reality the path of the Seine resembles a wave that enters Paris at the Quai Bercy, rises sharply northward, tapers slightly as it flows into separate streams around the islands, initiates its flat northernmost segment at the Place de la Concorde, then turns sharply in a great 60° bend at the Place d'Alma to flow out of the southwestern tip of the city. But in their drawings, 91.6 percent of the subjects understated the river's degree of curvature. Several subjects pulled it through the city as a straight line, and the typical subject represented the Seine as a gentle arc of slight but uniform curvature.

...of the river is made to resemble an arc of gentle

Perhaps, accurate full 3D is unnecessary?

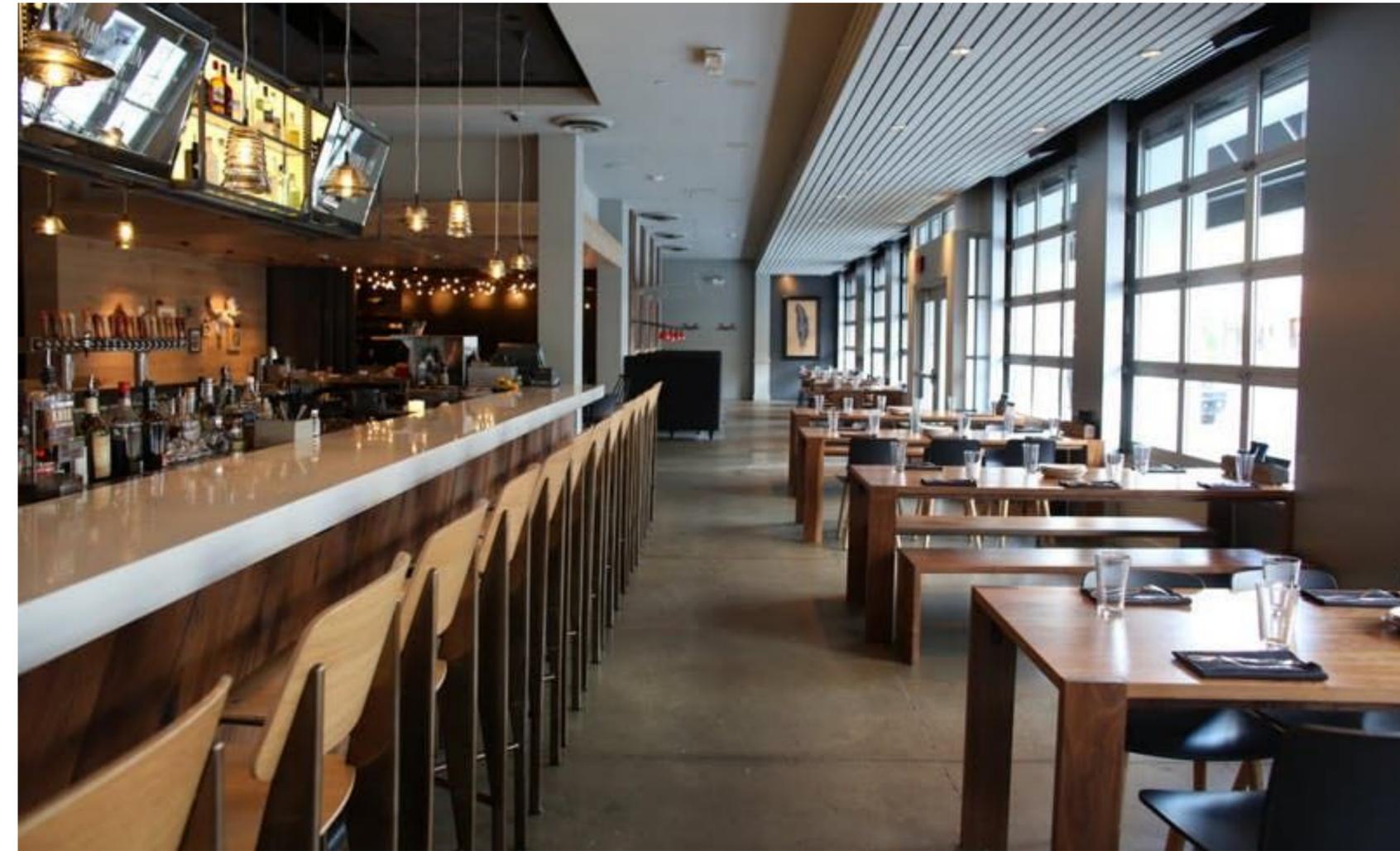
Lacks Semantics

Speculating about space not directly observed.



Lacks Semantics

Eg: Finding a bathroom in a new restaurant

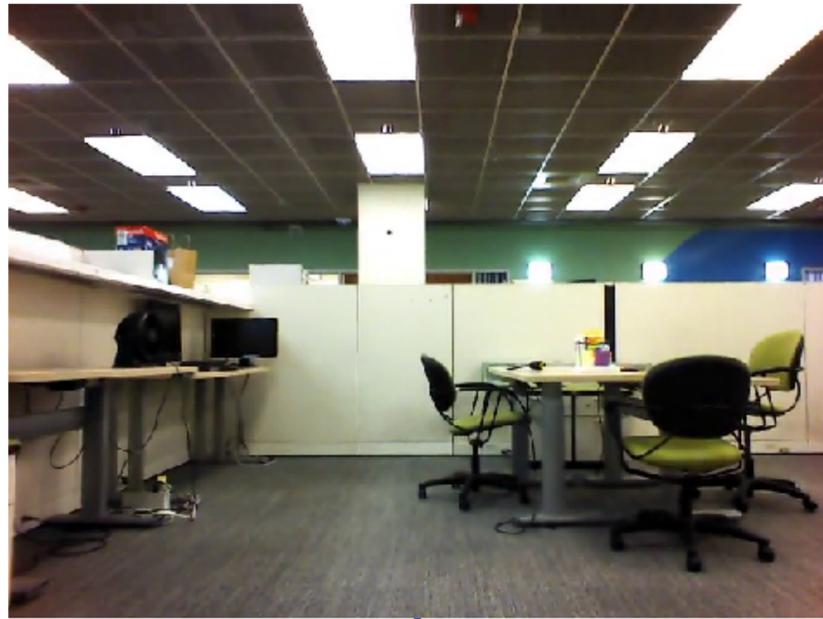


In this talk,

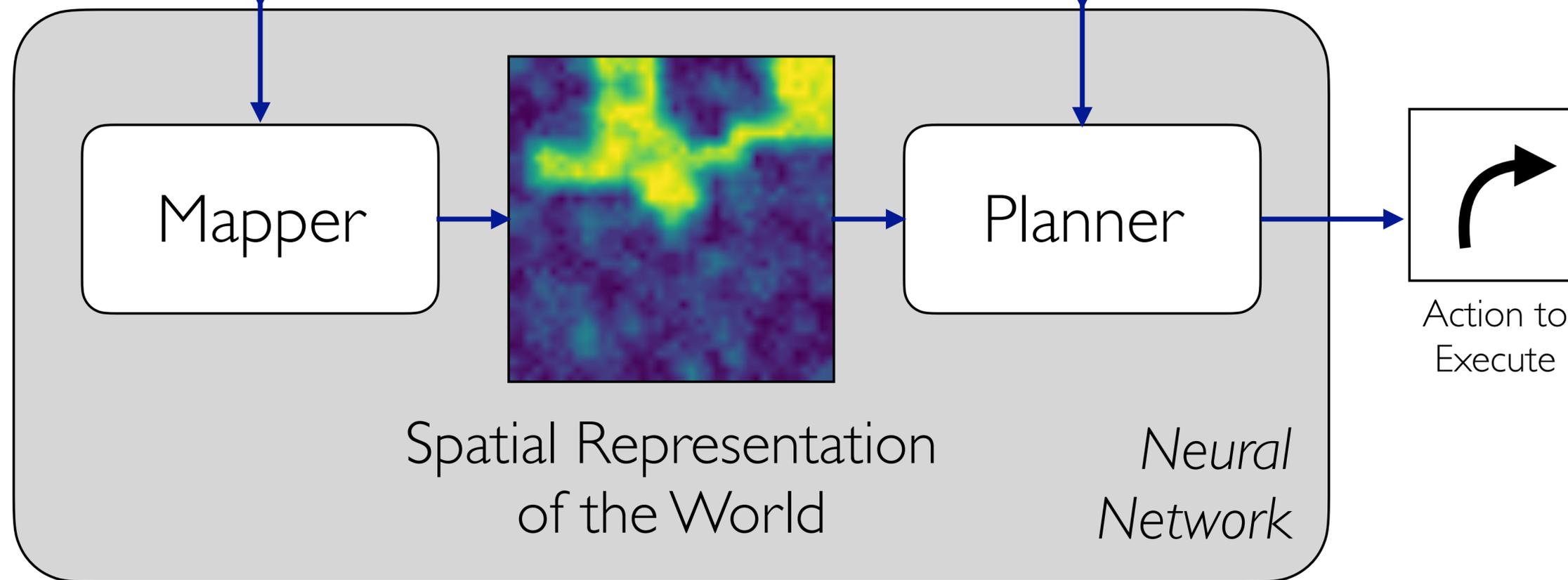
Representations for Places that Afford Navigation in Novel Environments

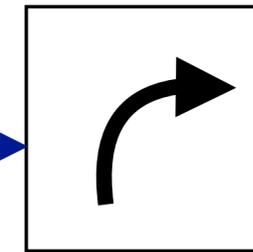
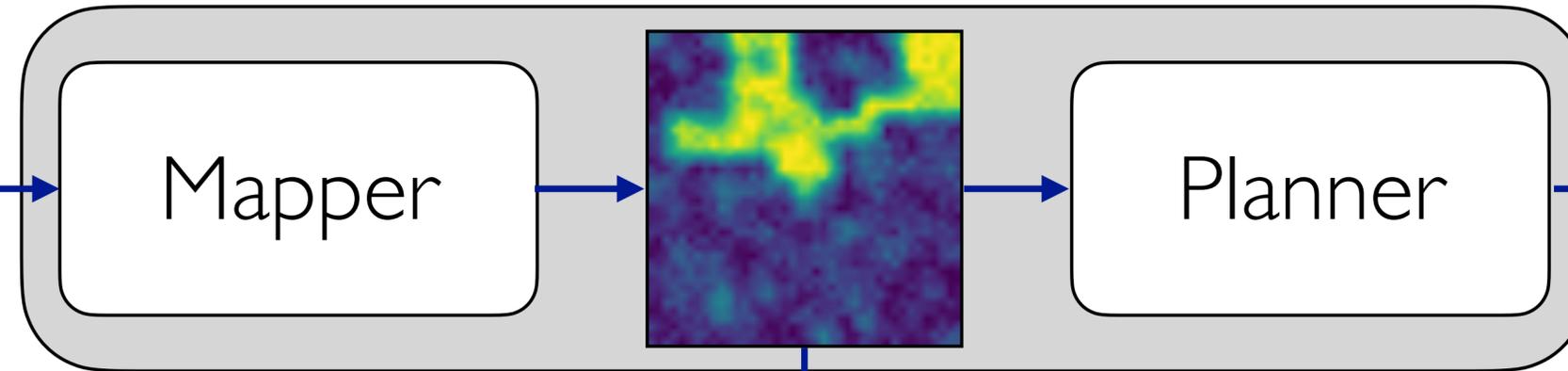
- *Augmenting metric representations with semantic reasoning*
- *Relaxing the need for metric representations*
- *Scaling-up training of such representations*

Operationalize insights from classical robotics into learning paradigms

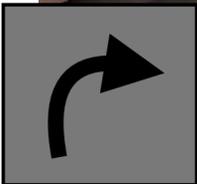
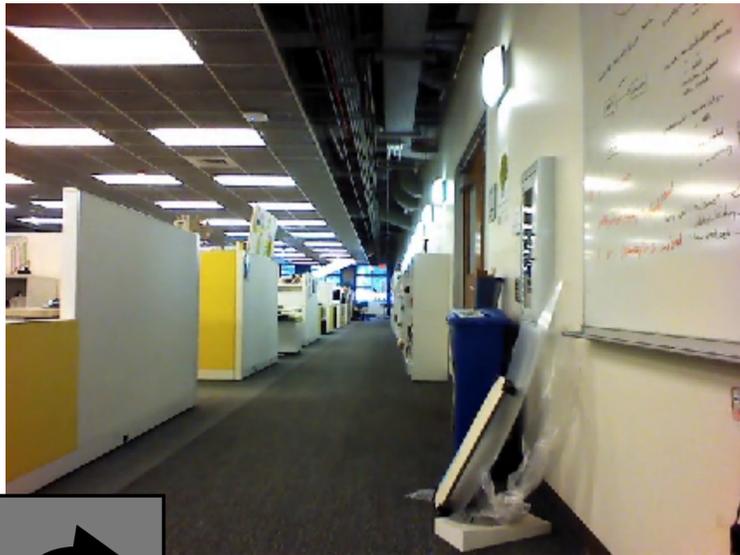


Goal (300, 400)

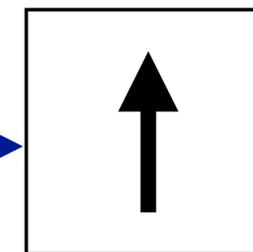
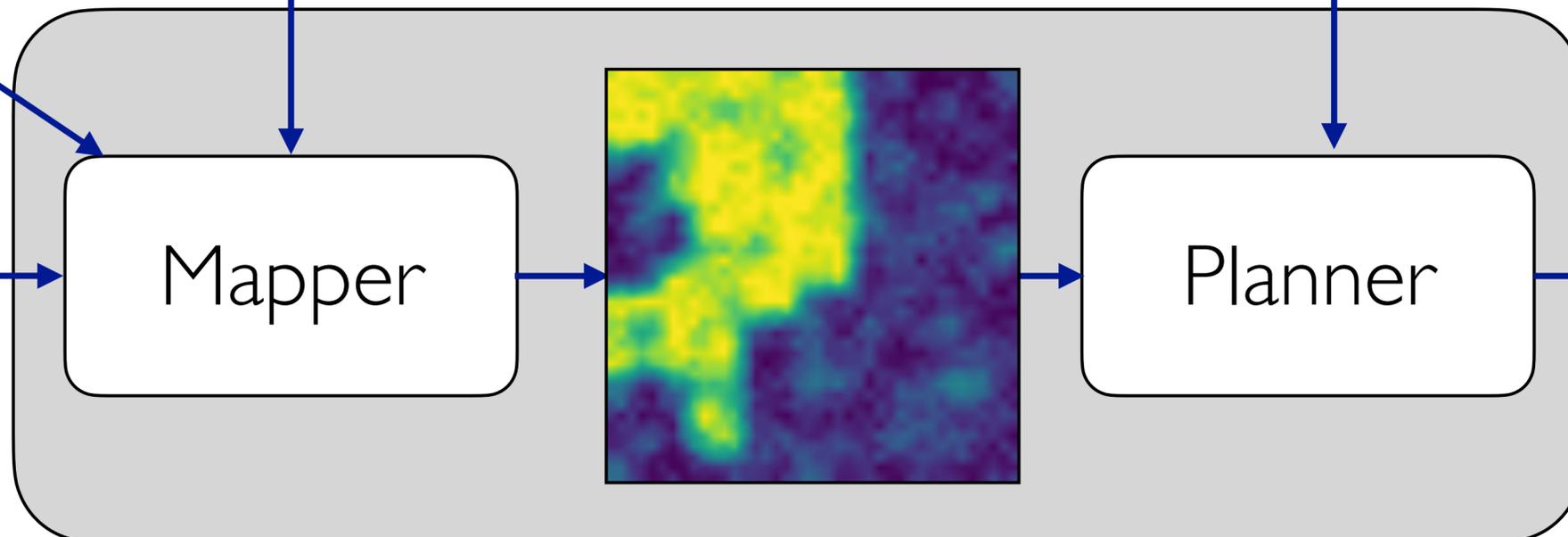




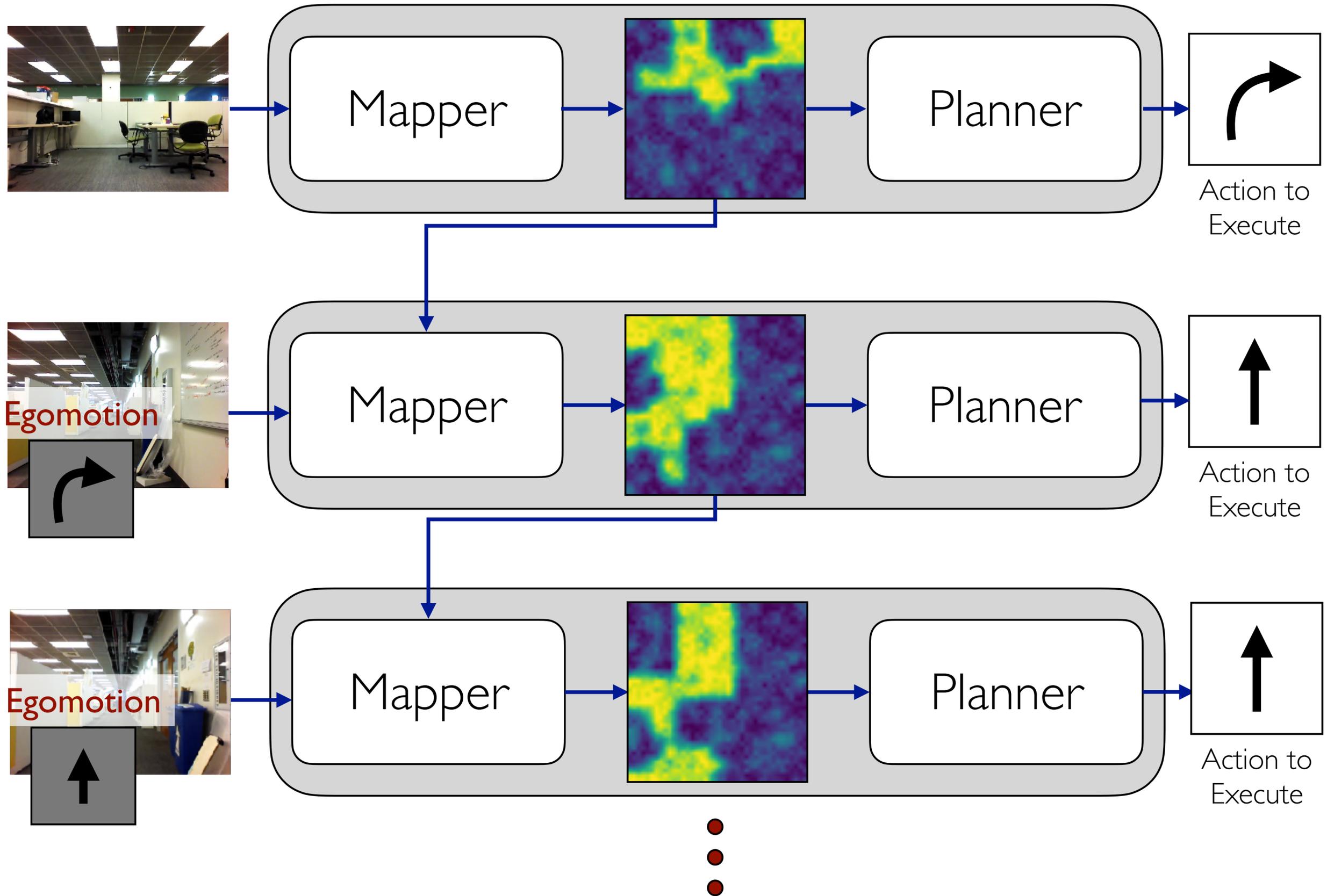
Action to Execute



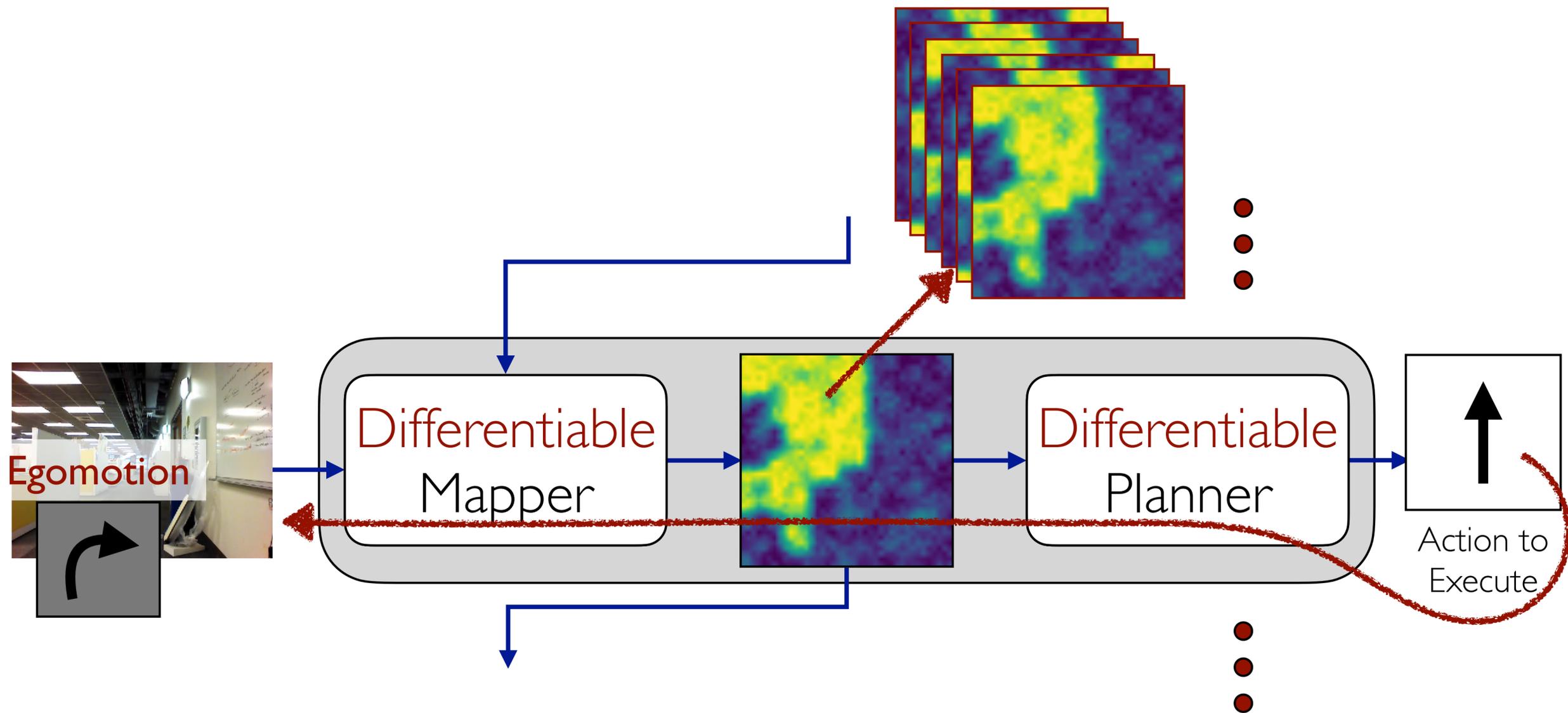
Egomotion



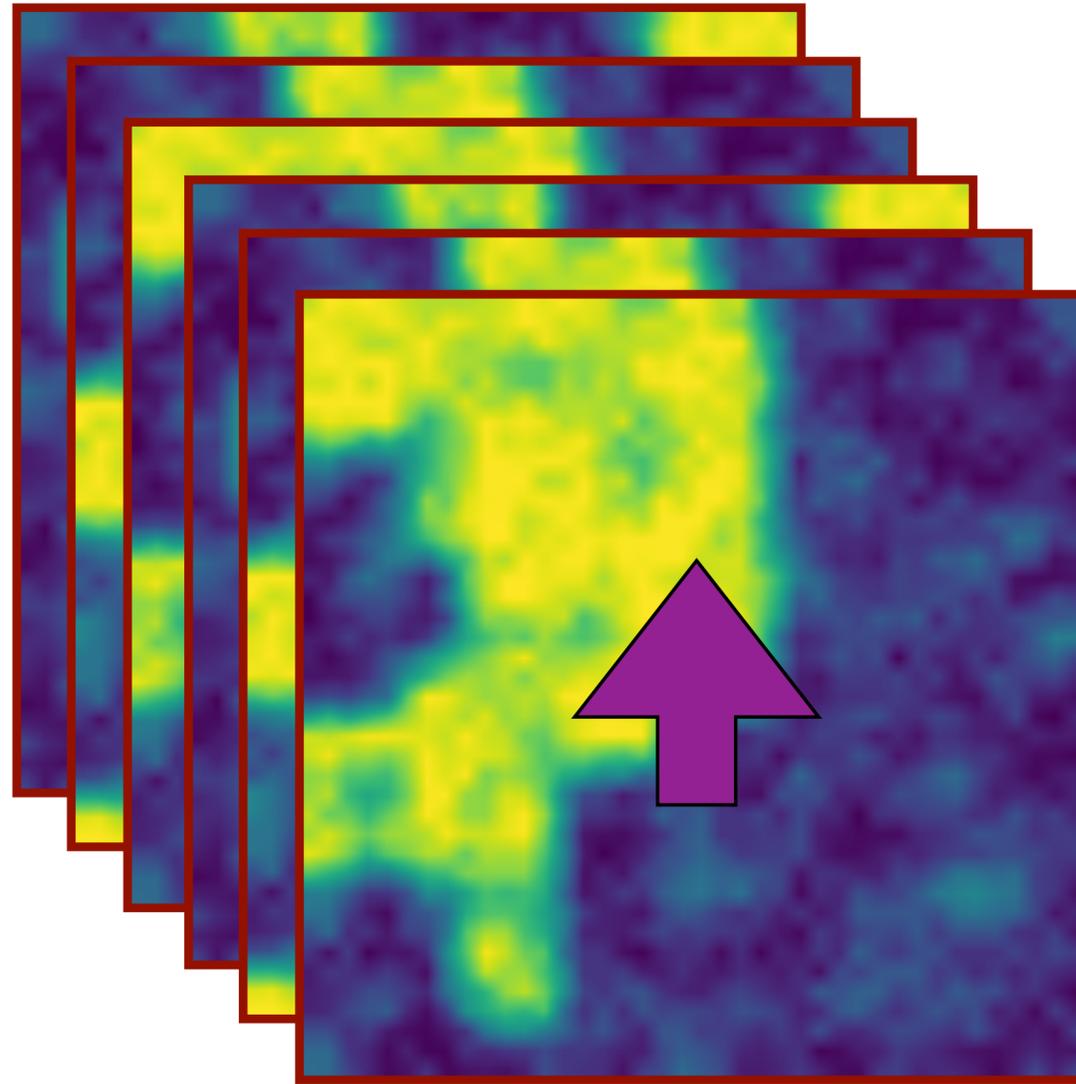
Action to Execute



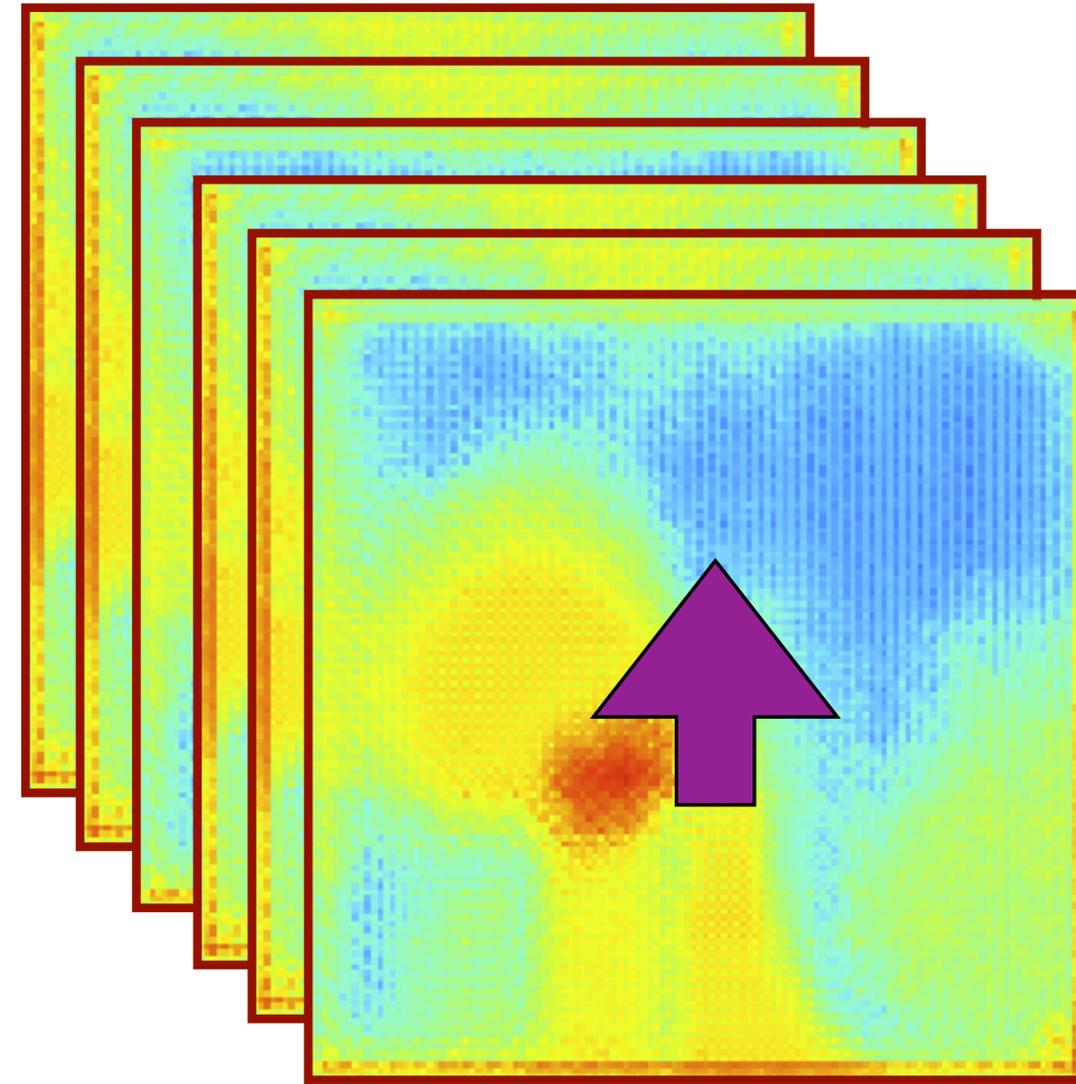
- Mapper and planner are *differentiable functions*
- Mapper and planner are *learned for end task*
- Hand-crafted obstacle maps to *task-driven semantic maps*



Spatial Representations



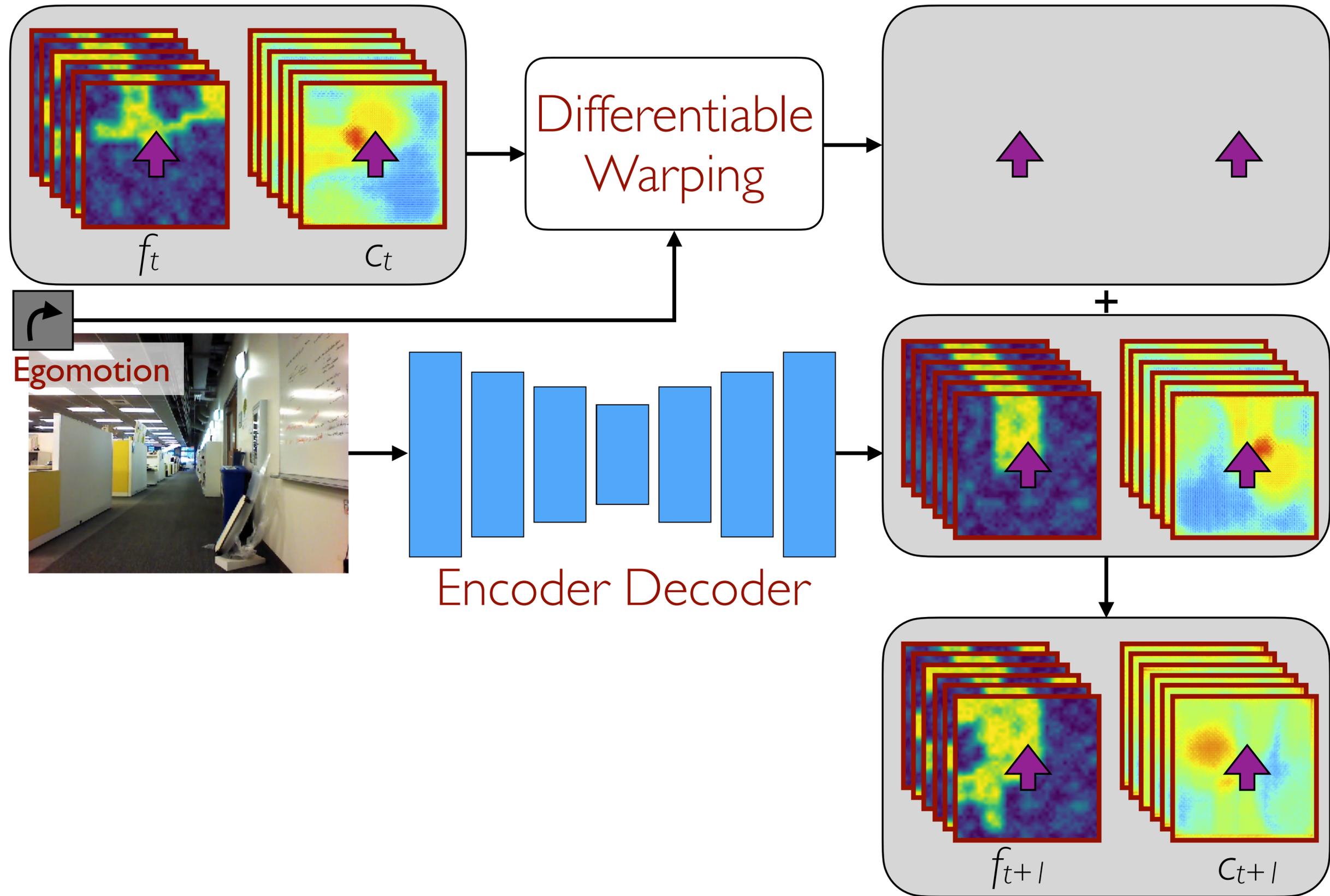
Feature f_t

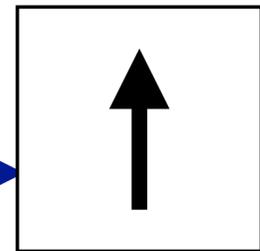
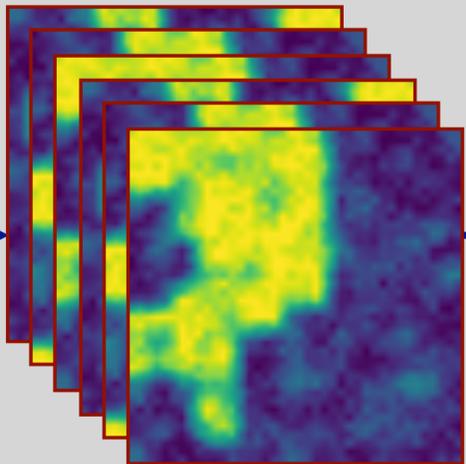
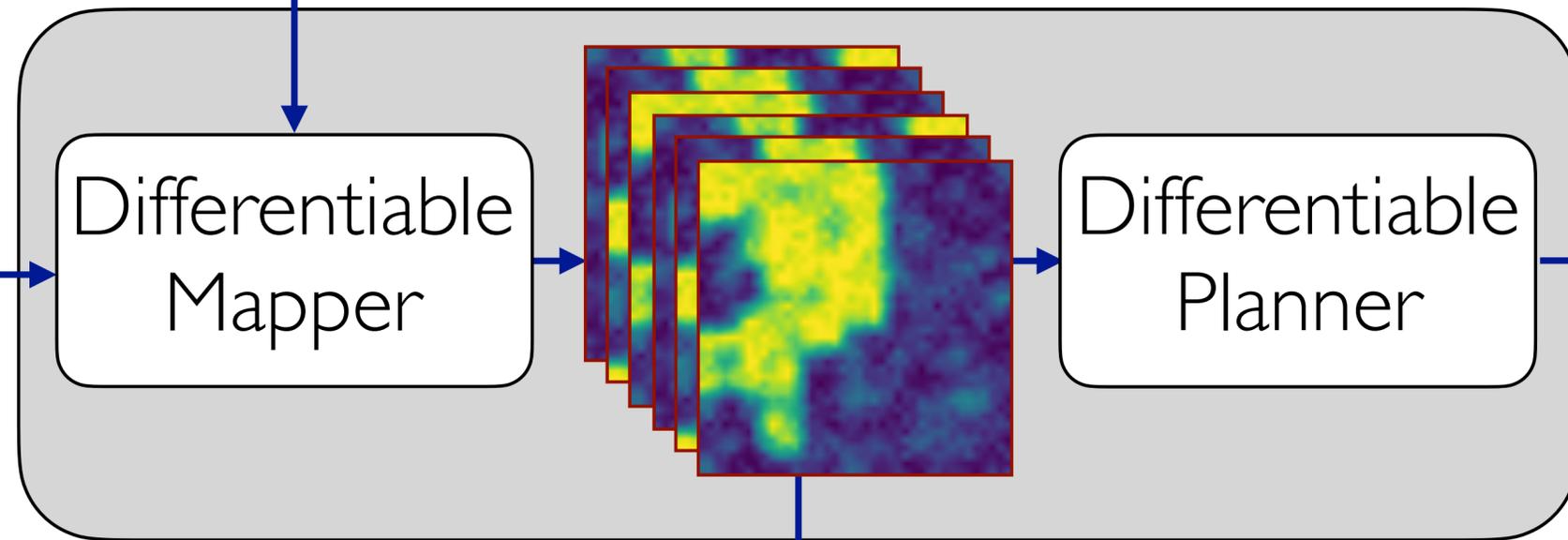


Confidence c_t

Egocentric Bird's Eye Coordinate Frame

Differentiable Mapper





Action to Execute



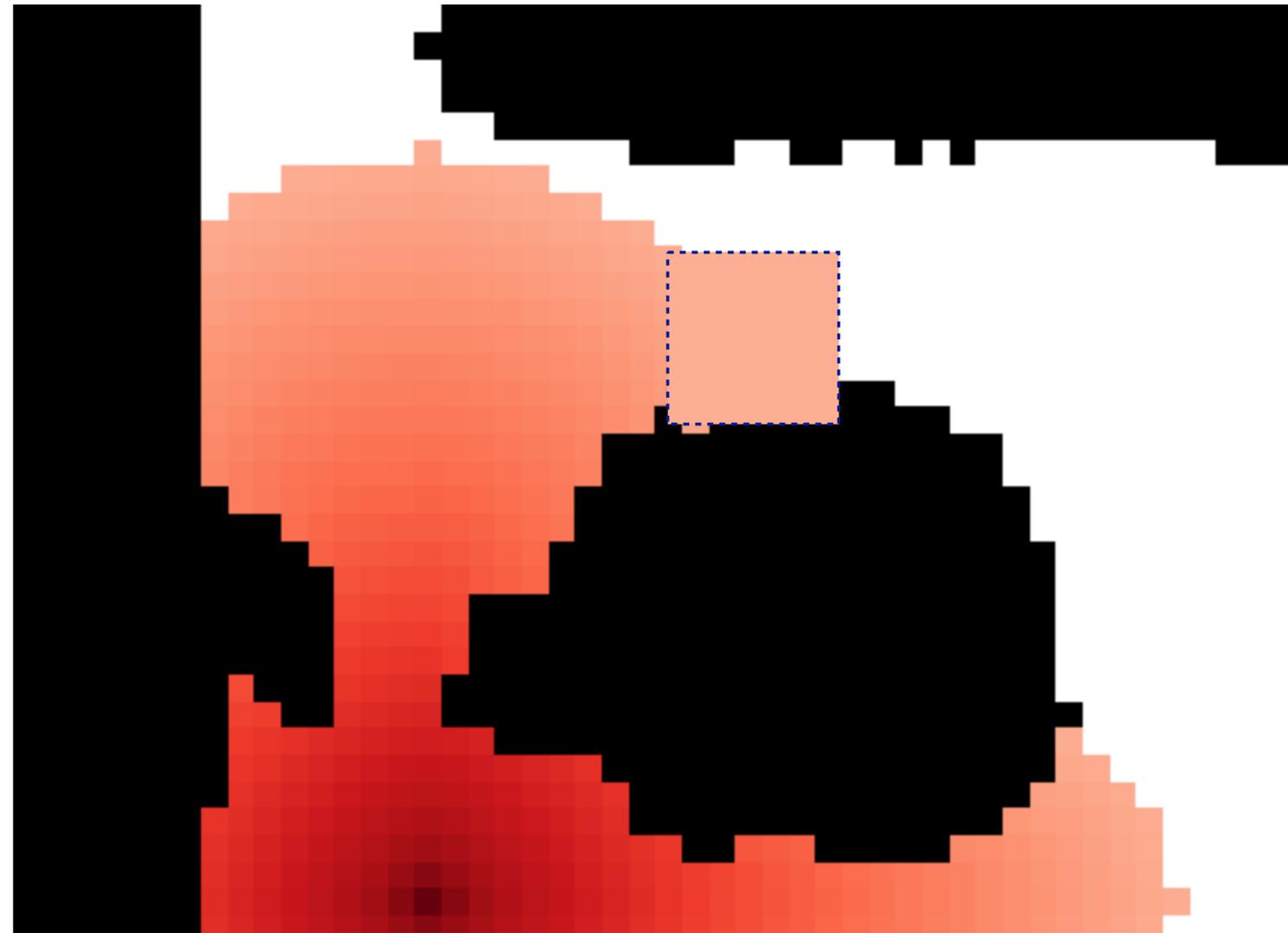
Differentiable Planner



Differentiable Planner



Differentiable Planner

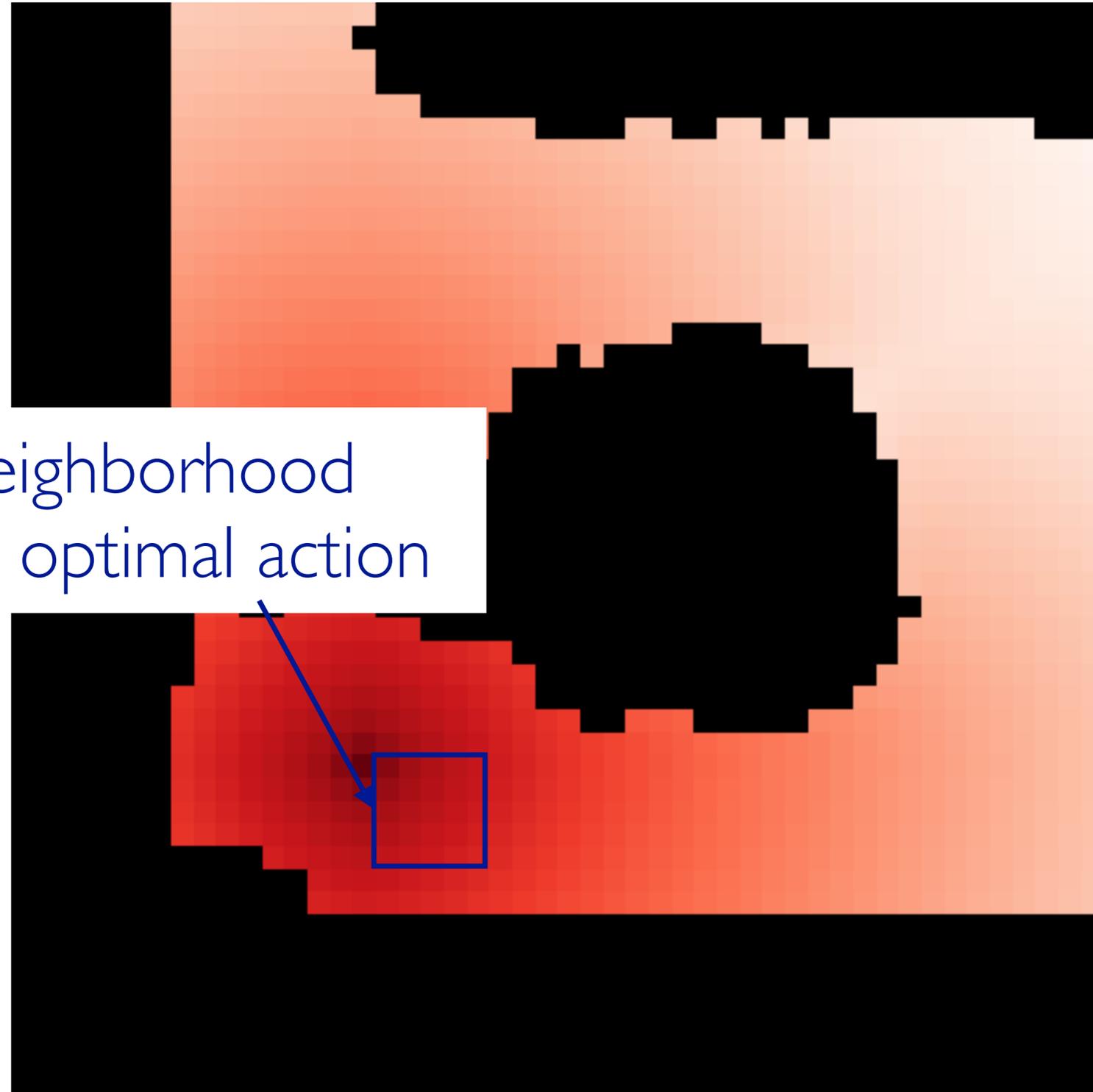


Local computation that can be done using
Convolutions and Channel-wise Max-Pooling.



Differentiable Planner

Local neighborhood
tells about optimal action



Policy Training

Simulator based on scans of *Real World Environments*



Simulate robot views and motion

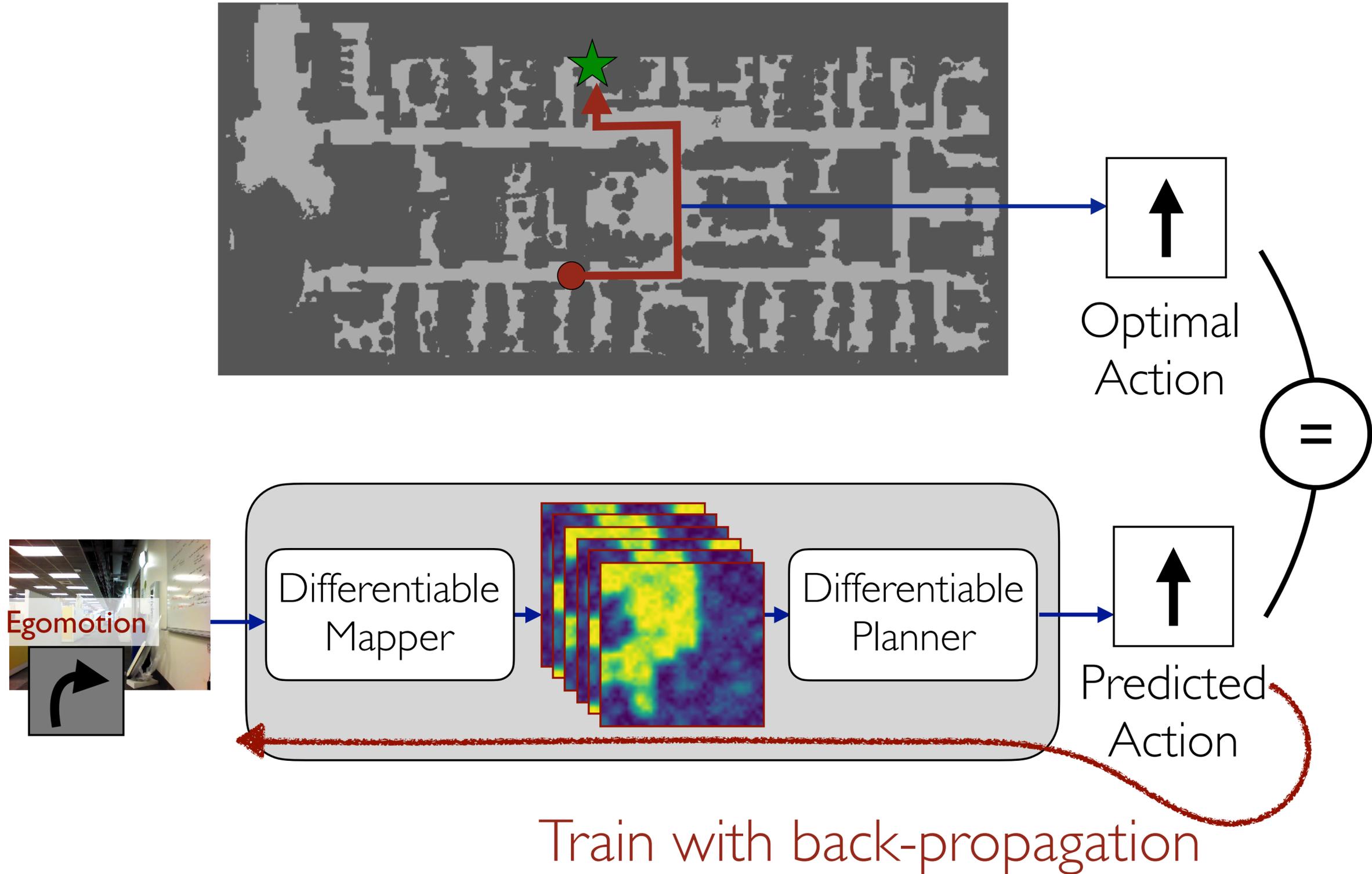


Compute ground truth traversability

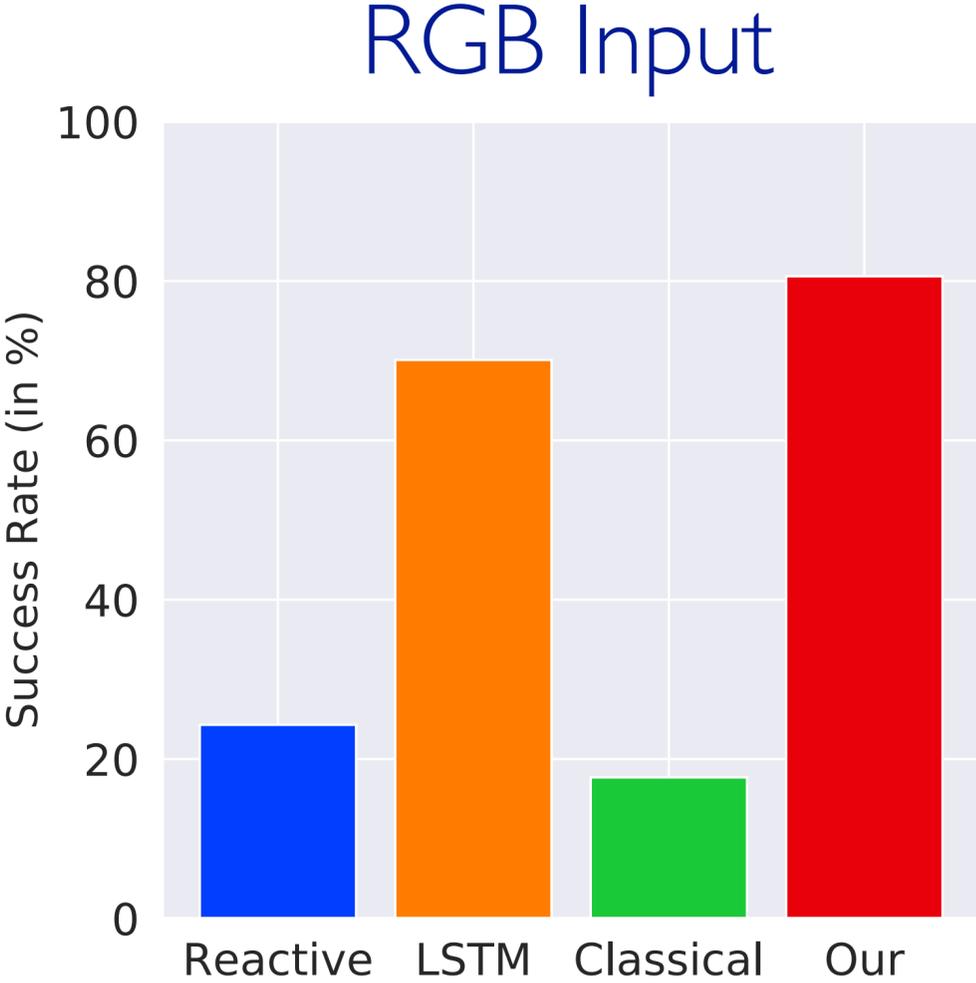
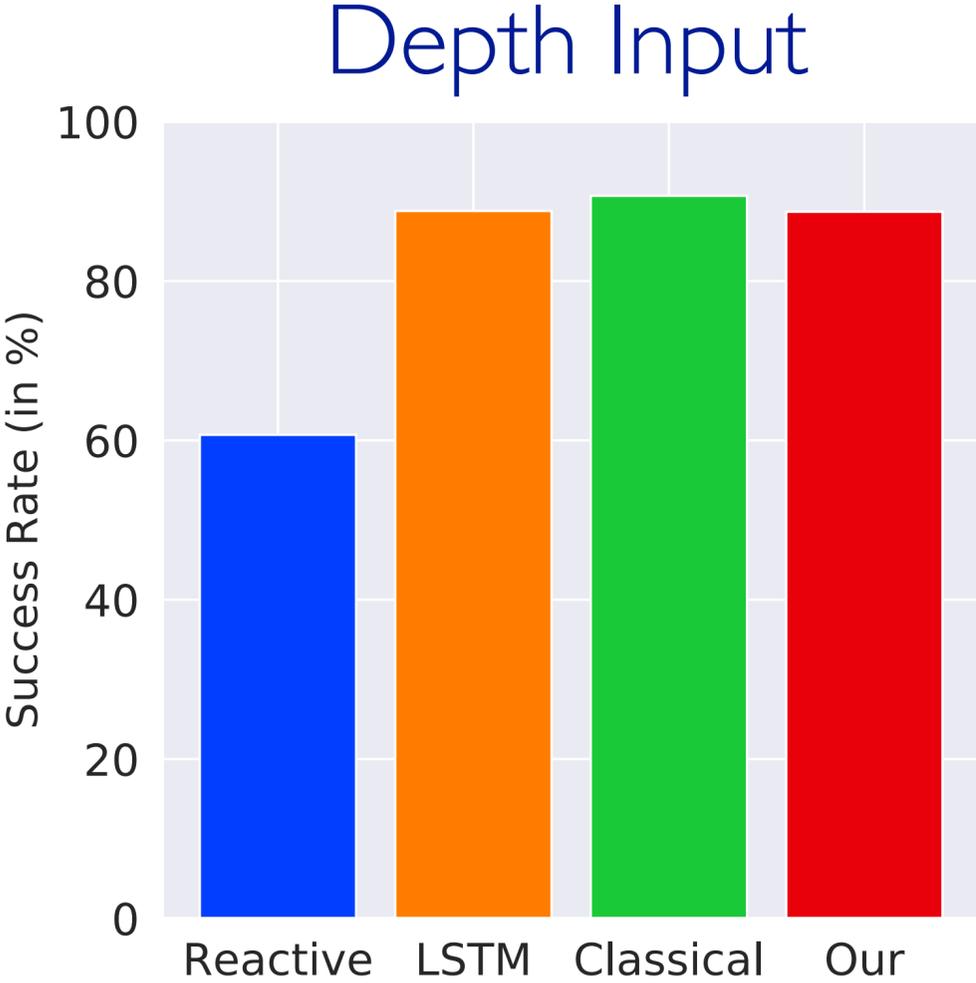
Gupta et al., CVPR 2017. Cognitive Mapping and Planning for Visual Navigation

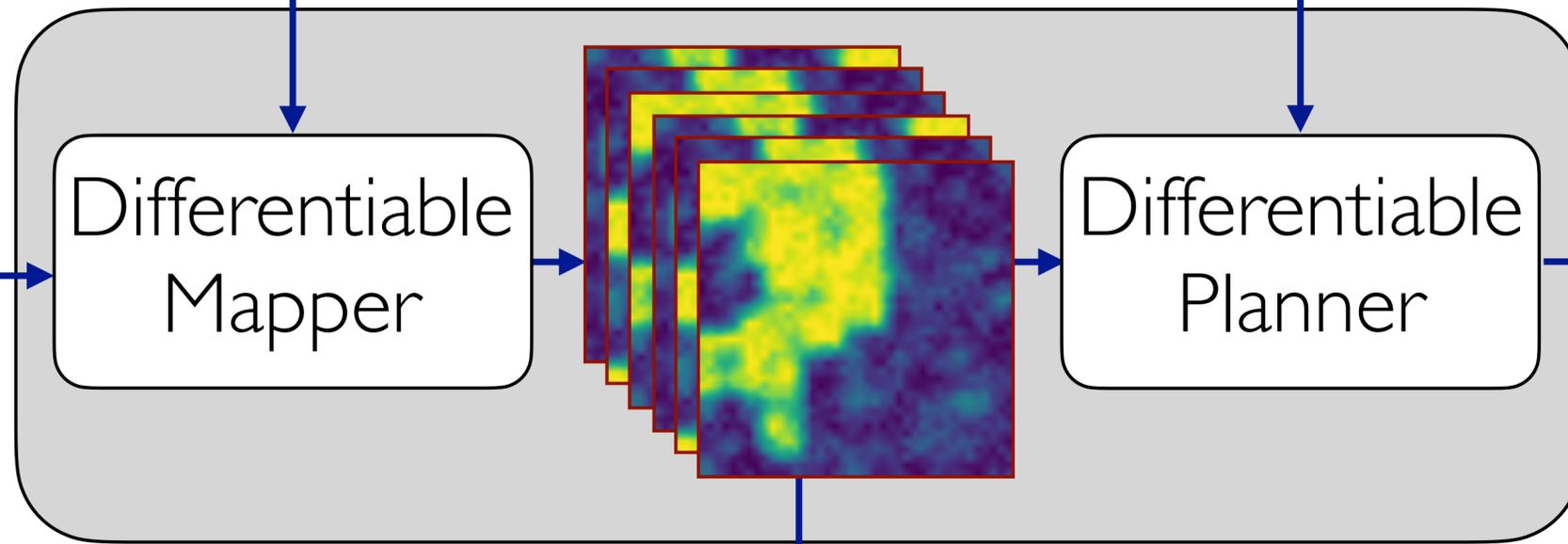
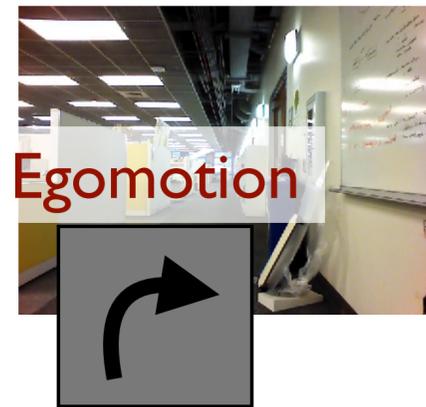
Armeni et al. CVPR 2016. 3D Semantic Parsing of Large-Scale Indoor Spaces

Policy Training by Expert Imitation



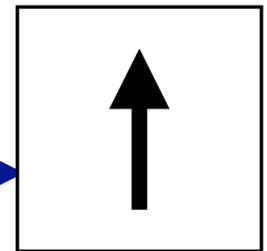
Results (Novel Env., Go To Relative Offset)





Goal (*chair*, table, door)

~~(70, 100, 100)~~

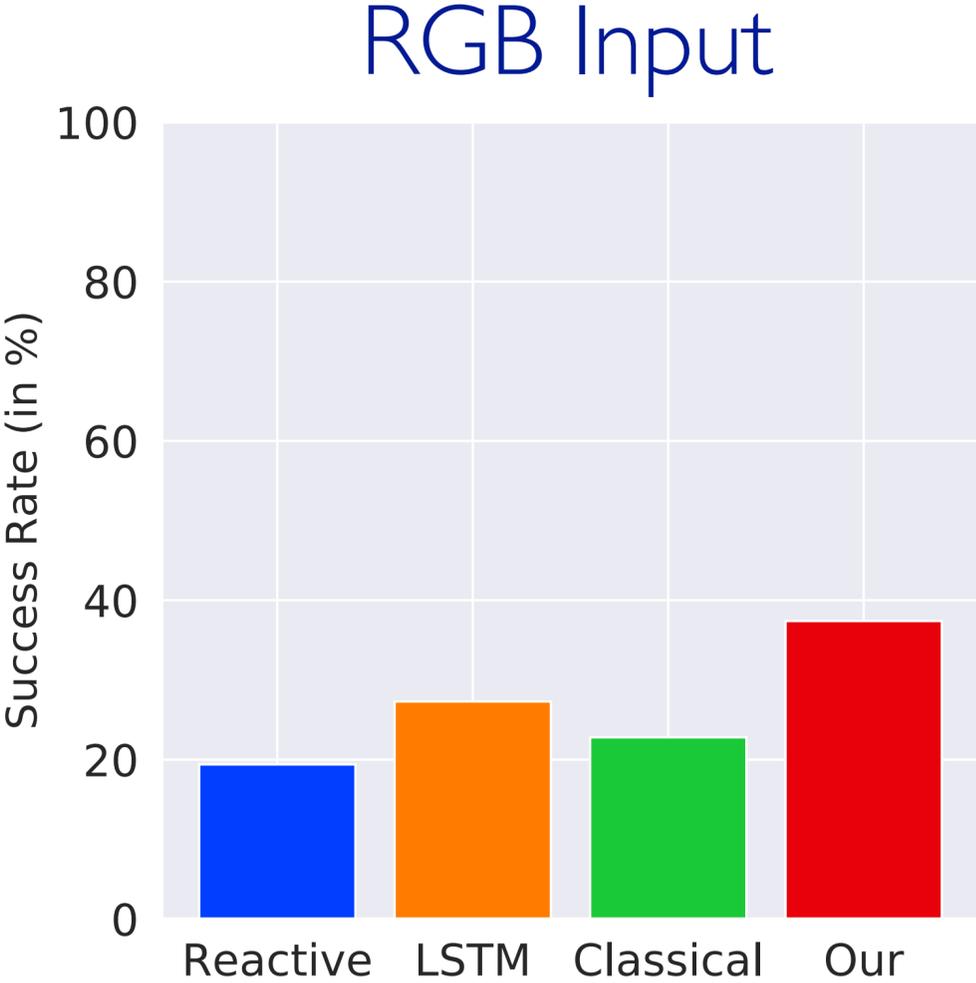
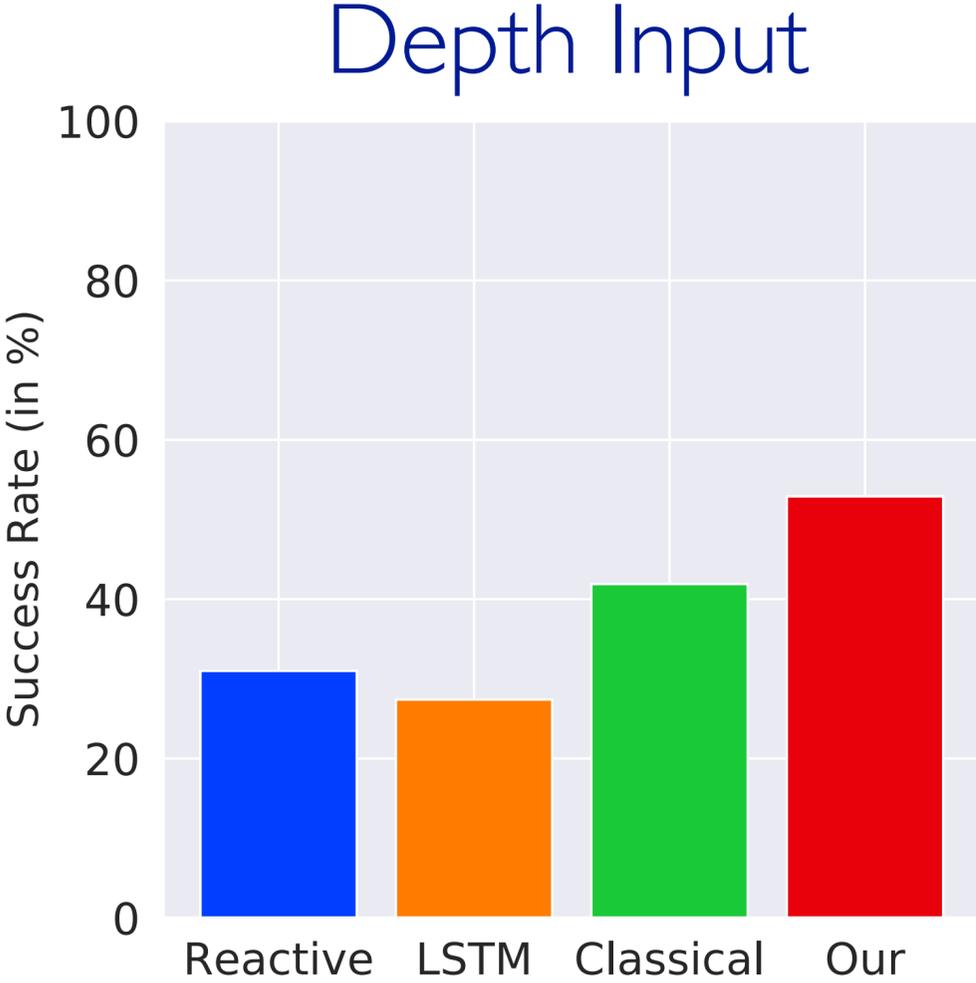


Action to Execute

Semantic Tasks (Go to a chair)

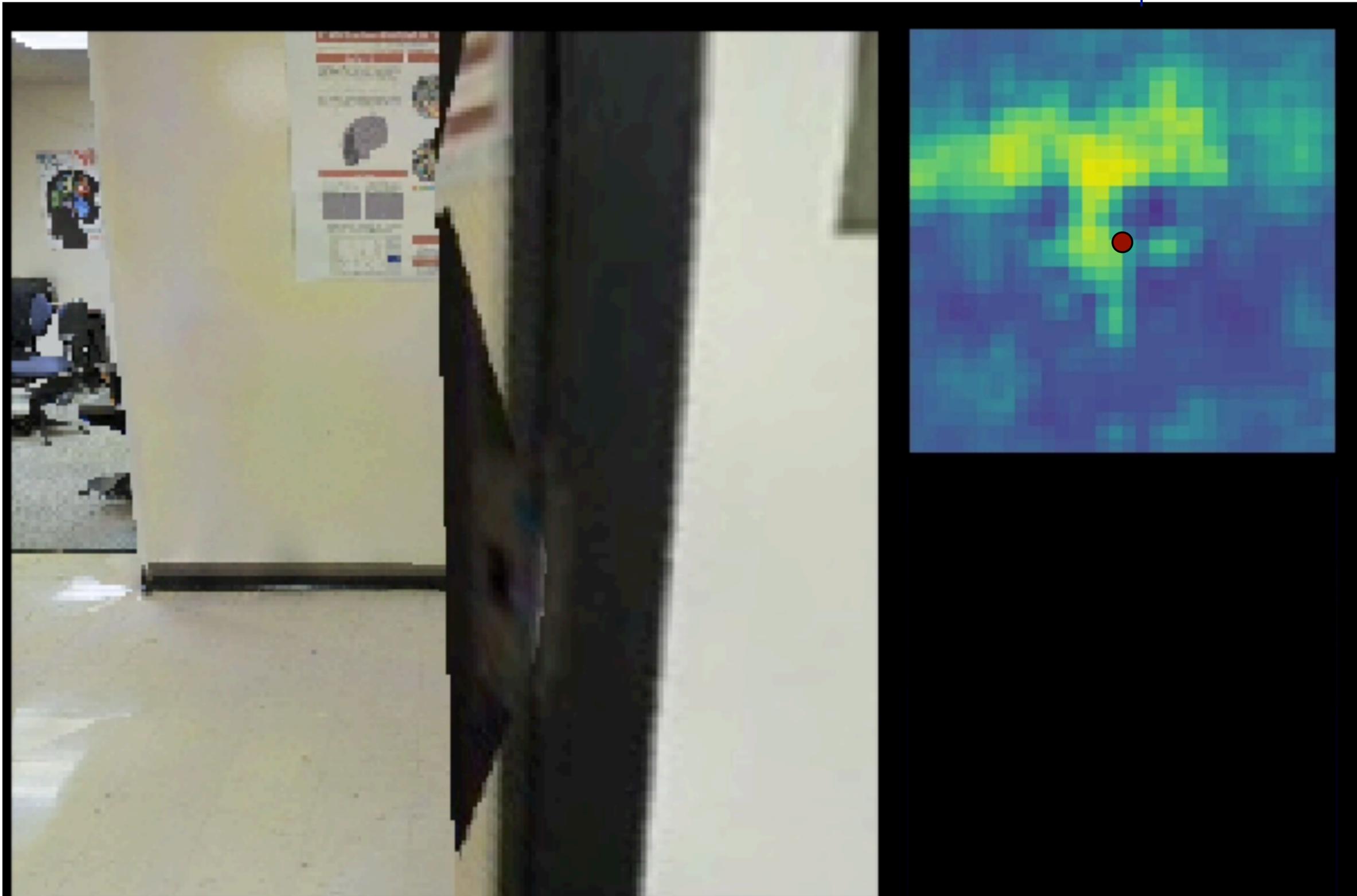
Successful Navigations
by CMP
(Semantic Task)

Results (Novel Env, Go To Object)



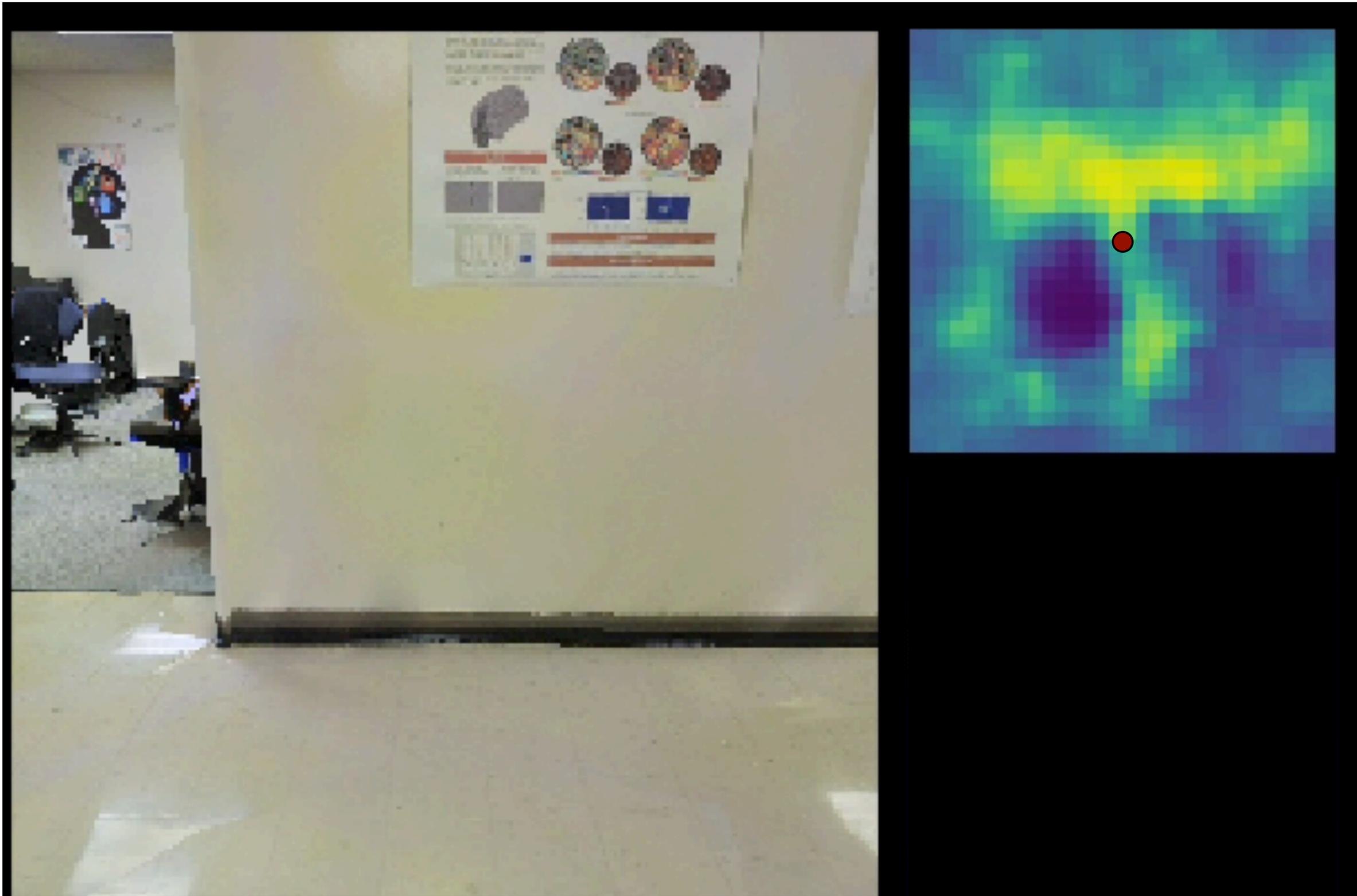
Agent can make predictions about its surroundings

Free Space



Agent can make predictions about its surroundings

Free Space

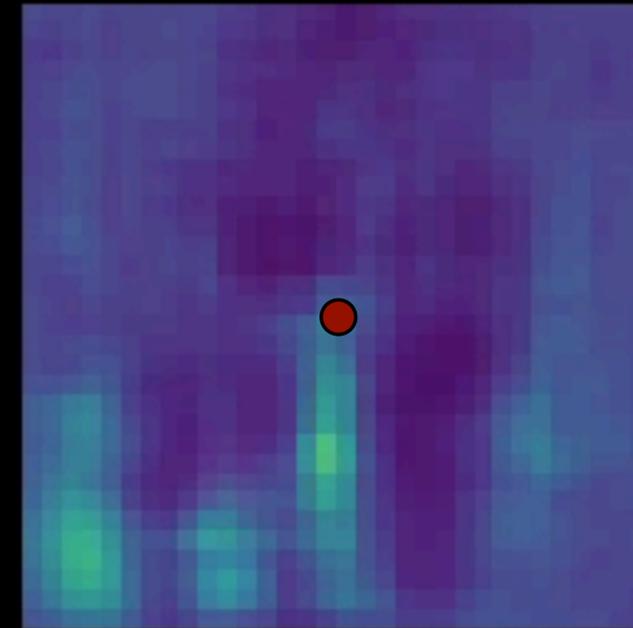
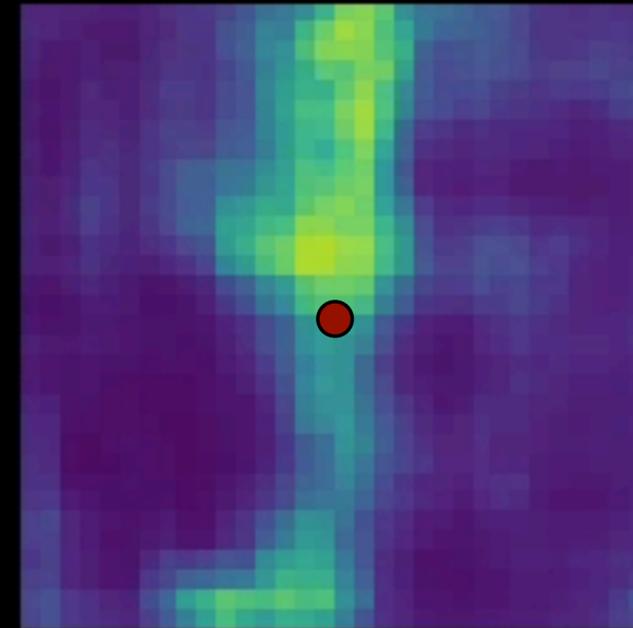
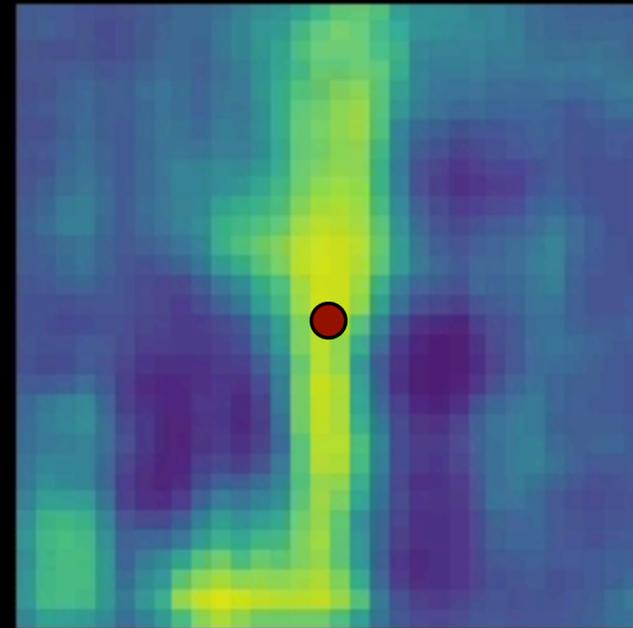
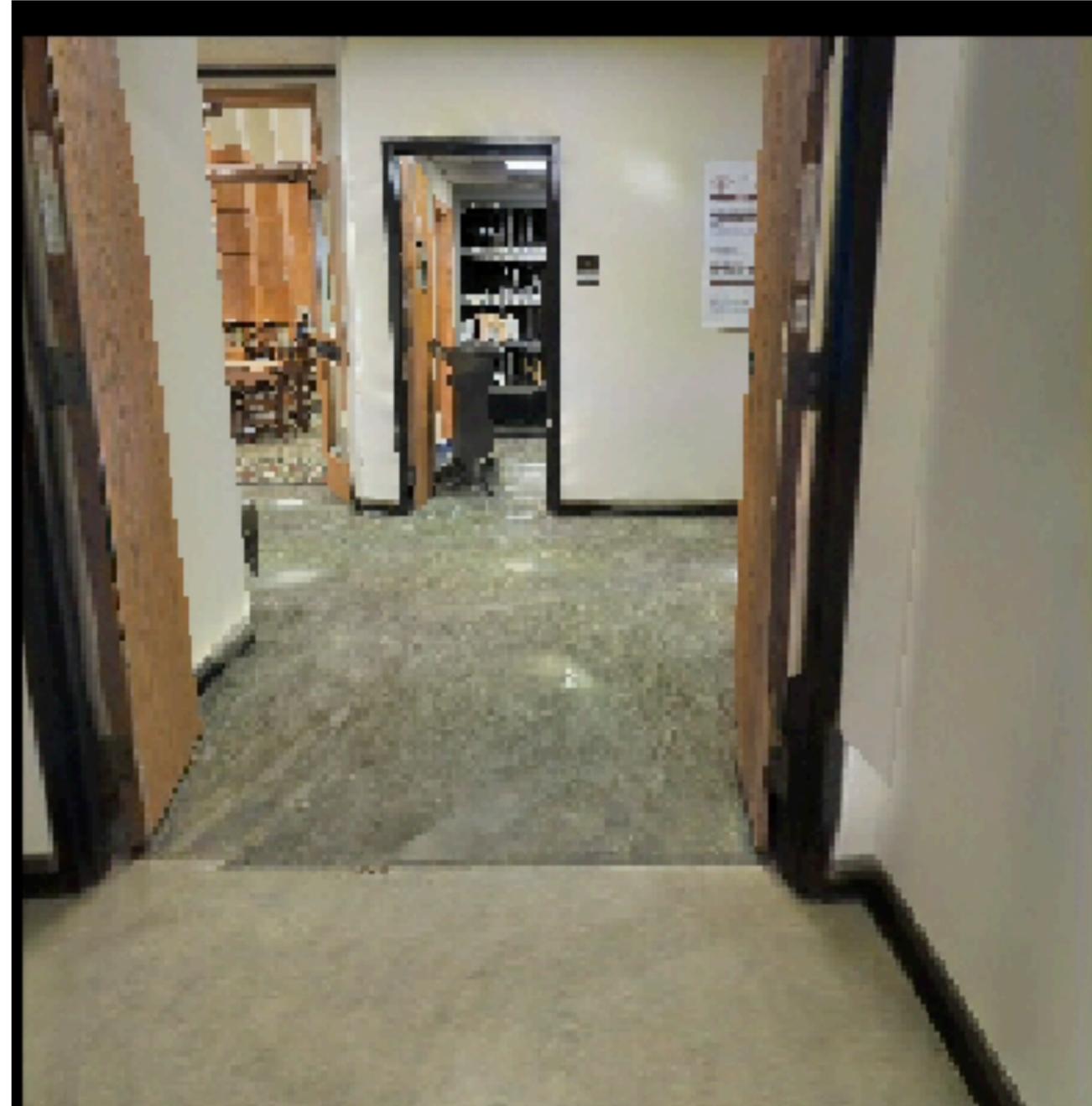


Agent can make predictions about its surroundings

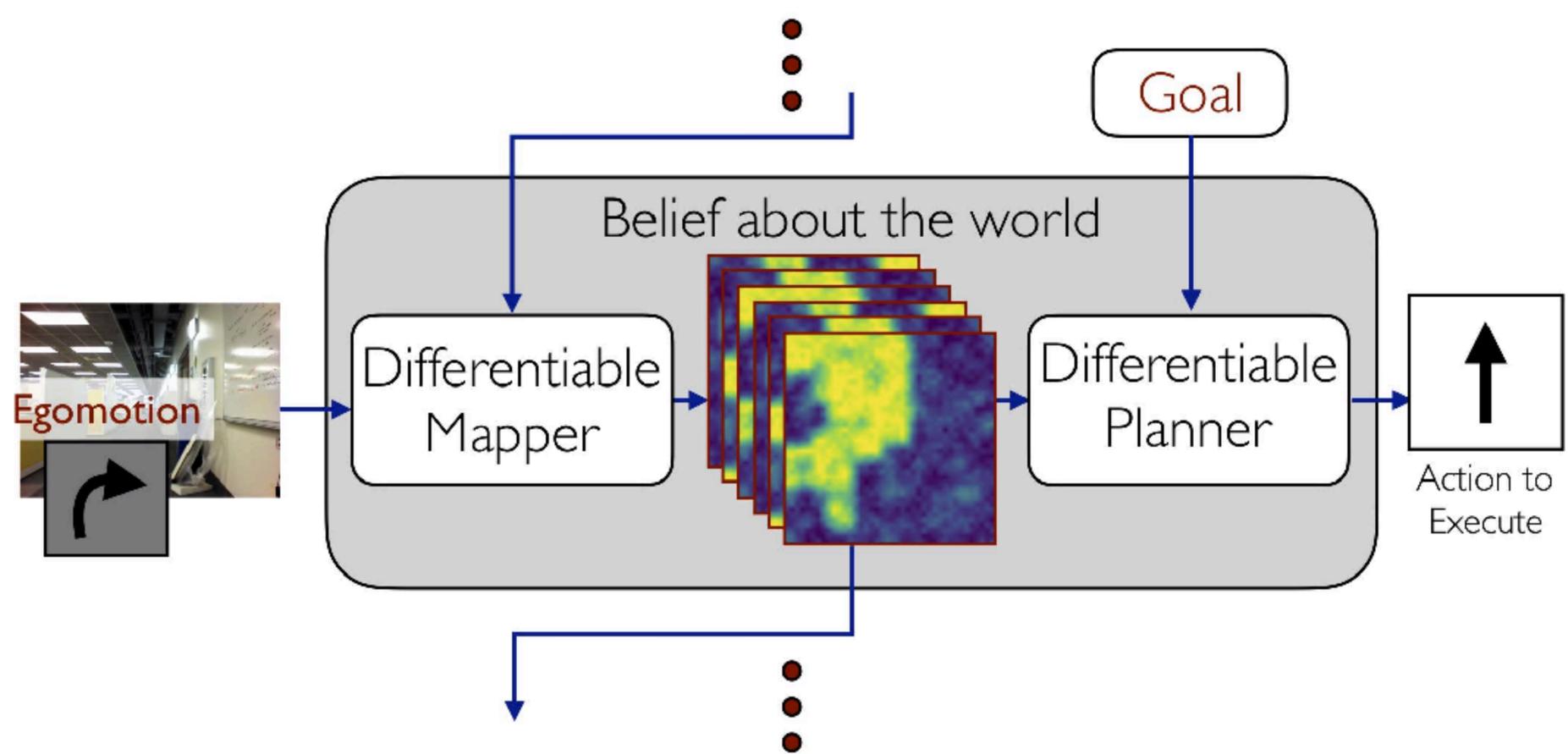
Free Space

Hallway

Room



Representation for Places



- Spatial reasoning
- Semantic reasoning
- Sensitive to pose error
- Interactive training (DAgger easier than RL, but still)
- Long training horizons

Can we relax the need for spatially consistent global maps?

Reaching Image Goals



Image Goal Task

- Agent observations are panoramic images
- Take actions to navigate to the goal location
- Take the 'stop' action at the goal location
- *Actuation noise, robot does not precisely know how much it has moved*

Source Image, I_s



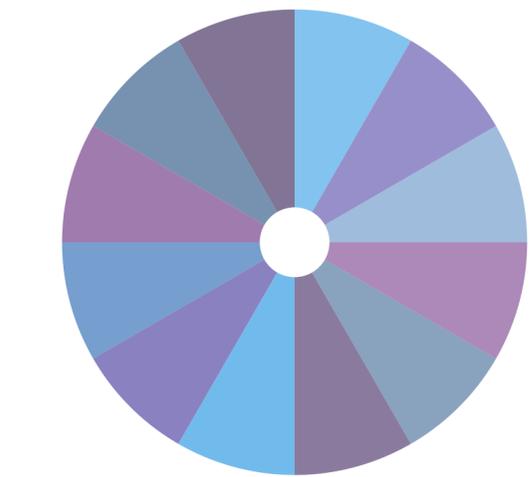
navigation
actions



Target Image, I_g

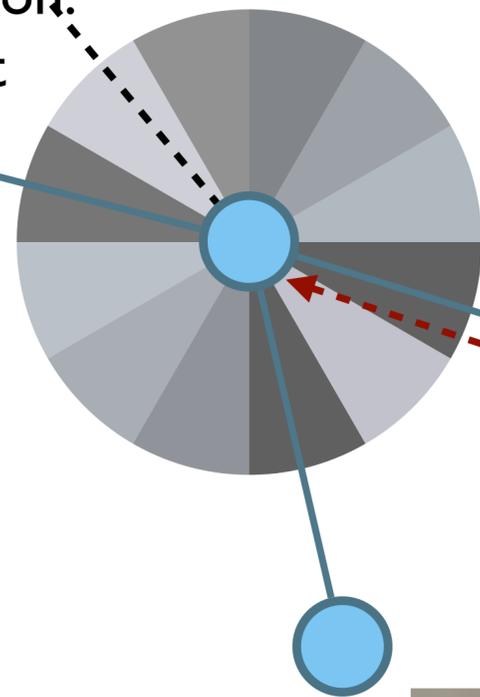


Representation

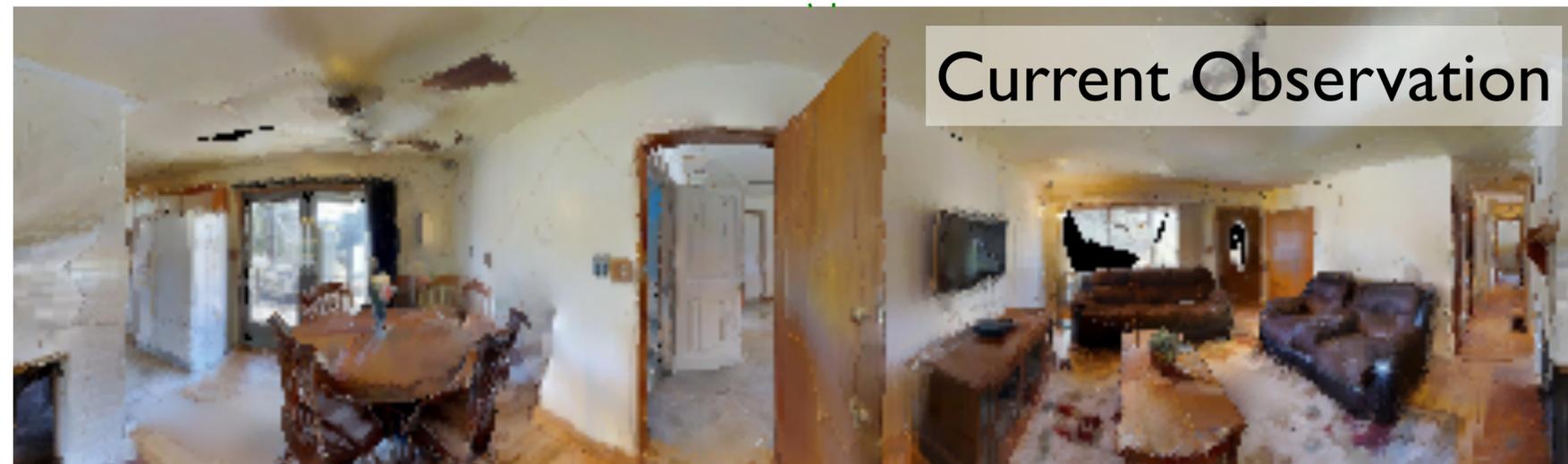
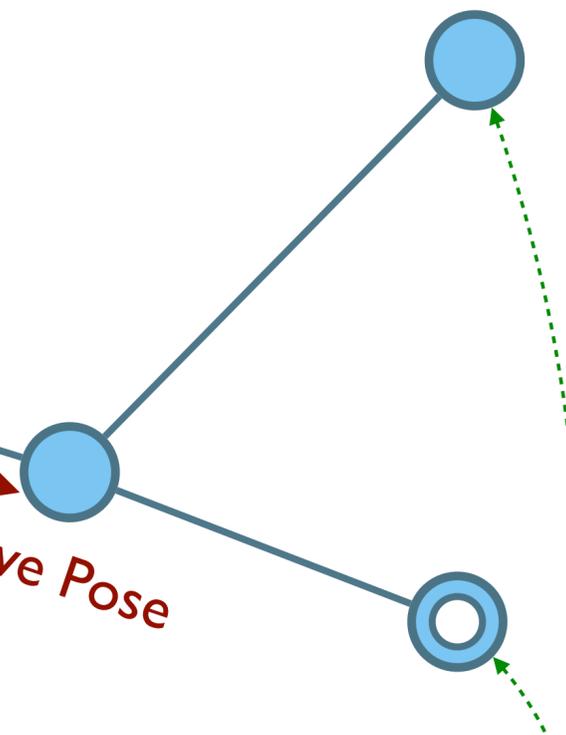


$F_G = \text{Geometric Prediction: Free space in different directions}$

"Ghost Nodes" 0.8



$F_L = \text{Localization}$



4 Functions

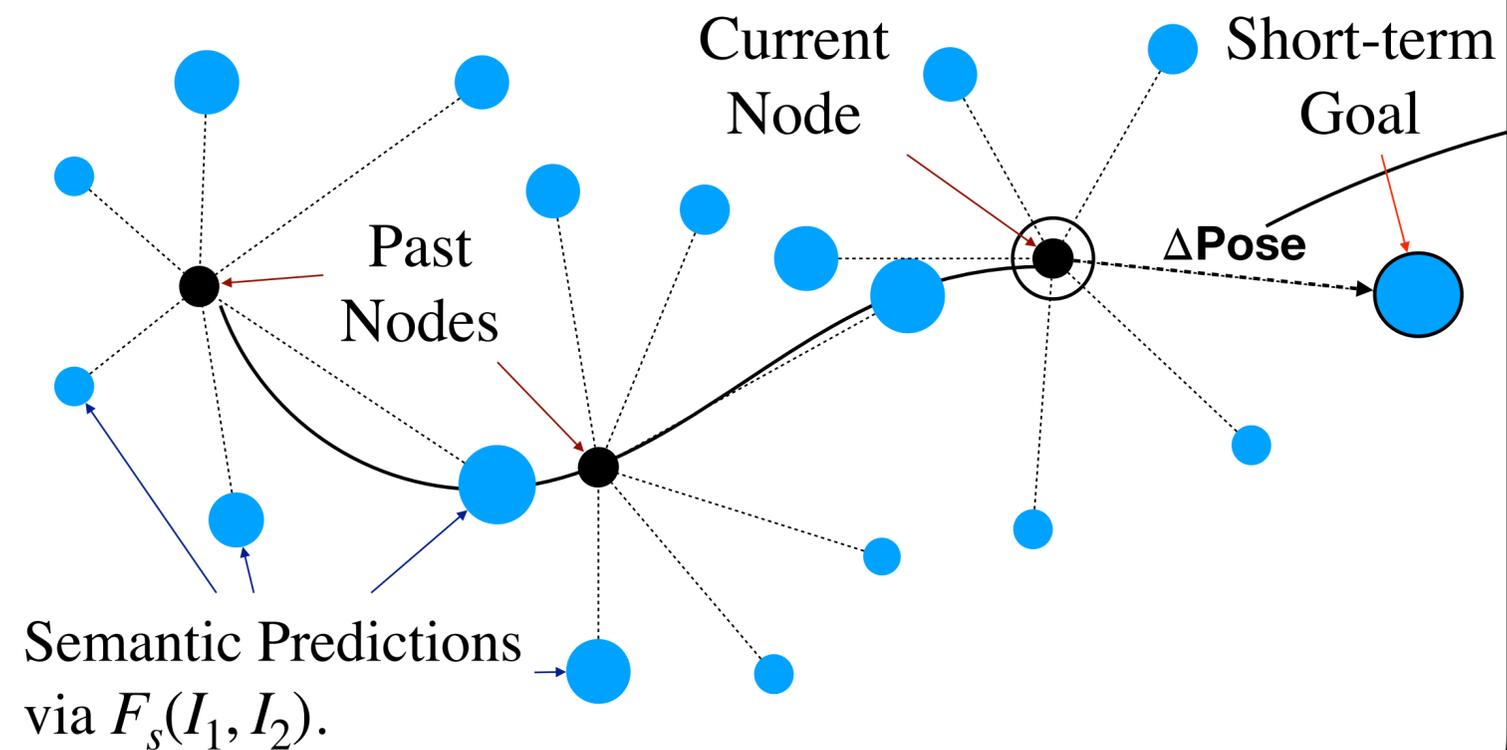
- $F_G(I_1)$ = Geometric Prediction: Free directions
- $F_S(I_1, I_2)$ = Semantic Prediction: Closeness to target
- $F_R(I_1, I_2)$ = Relative Pose
- $F_L(I_1, I_2)$ = Localization

Using the Representation

Hierarchical Policy

High-Level Policy

- Decides where to go next and emits short-term goal
- Builds a topological map that keeps track of visited nodes, ghost nodes and values predicted by $F_s(I_1, I_2)$ for different directions



Low-Level Policy

- Executes actions to achieve short-term goal
- Option 1: predict local occupancy, plan paths



Occupancy Map

FMM Cost Map

- Option 2: learn a low-level controller

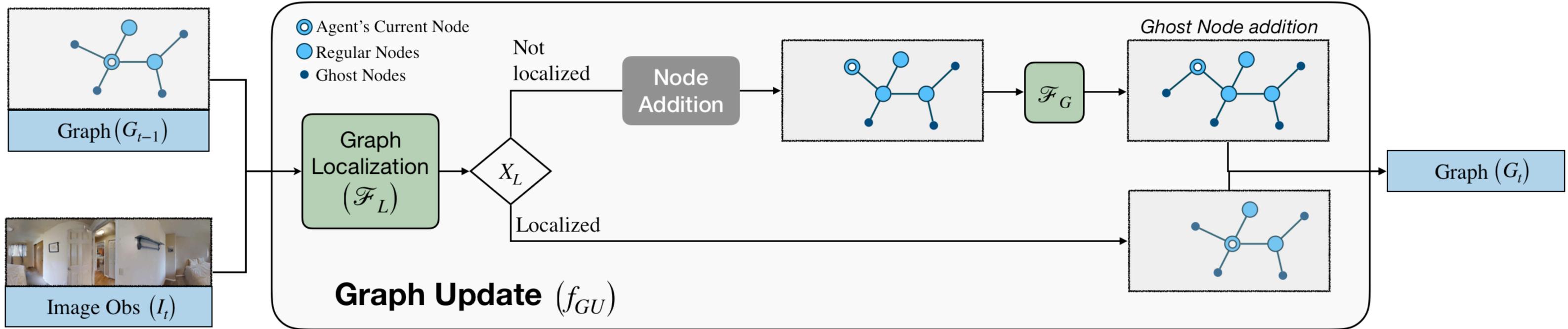
Forward

Left

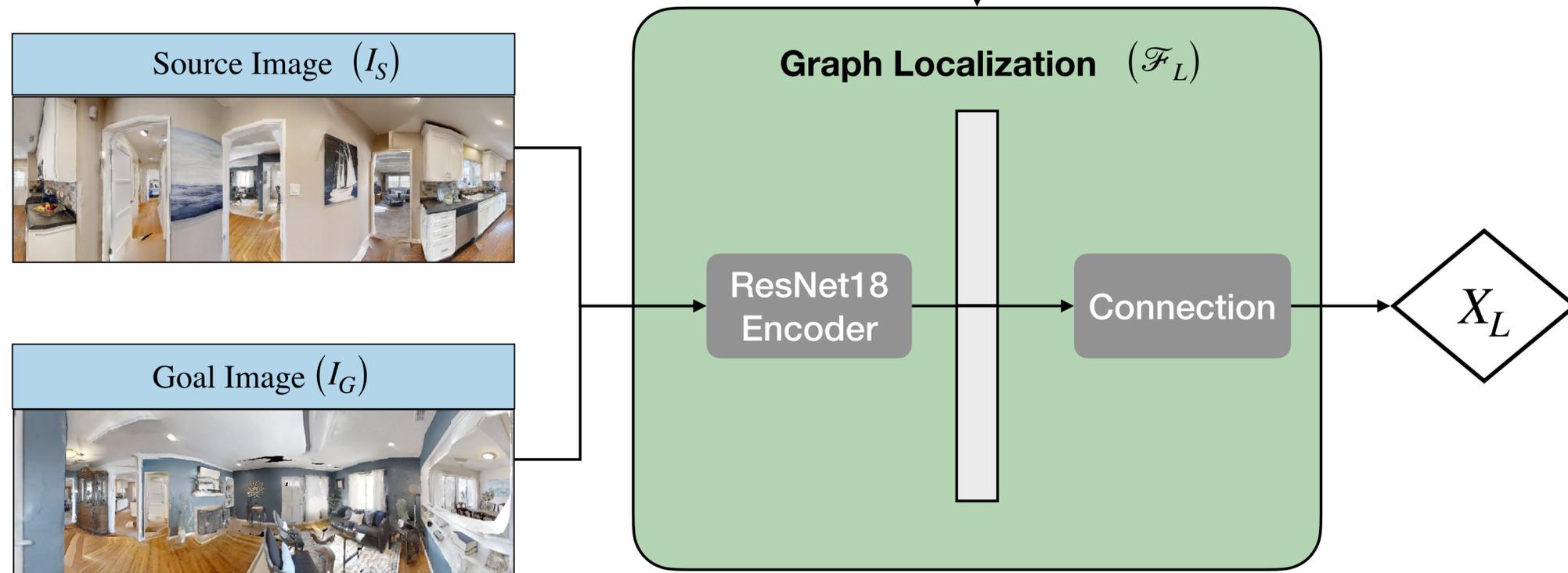
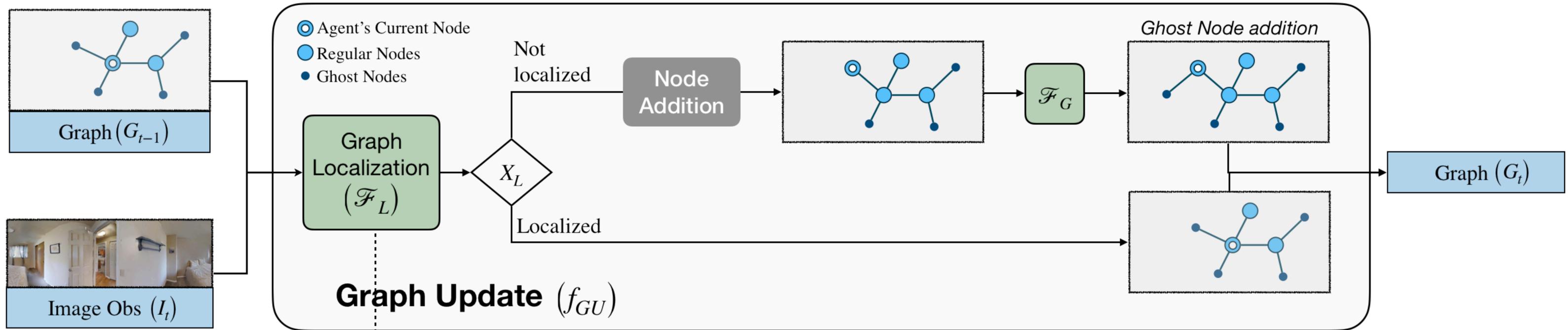
Right

Stop

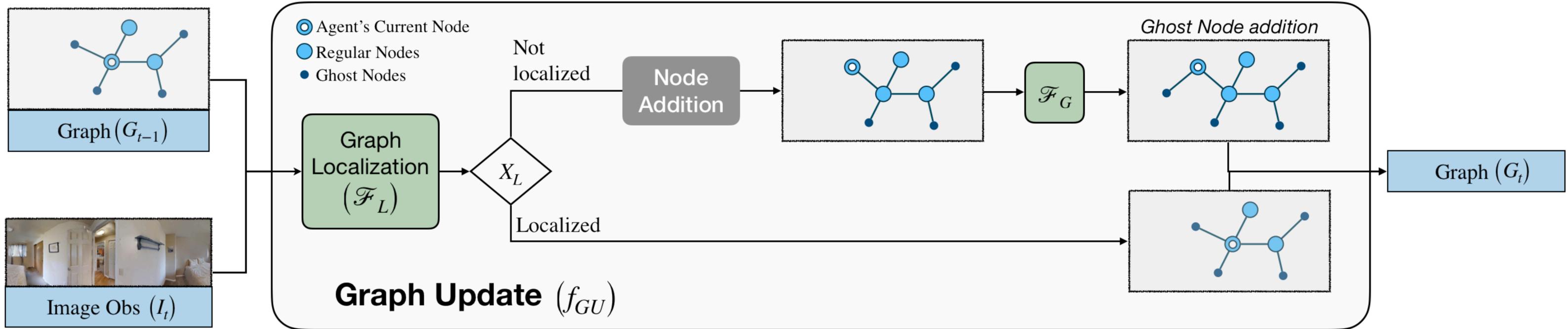
Building the Representation



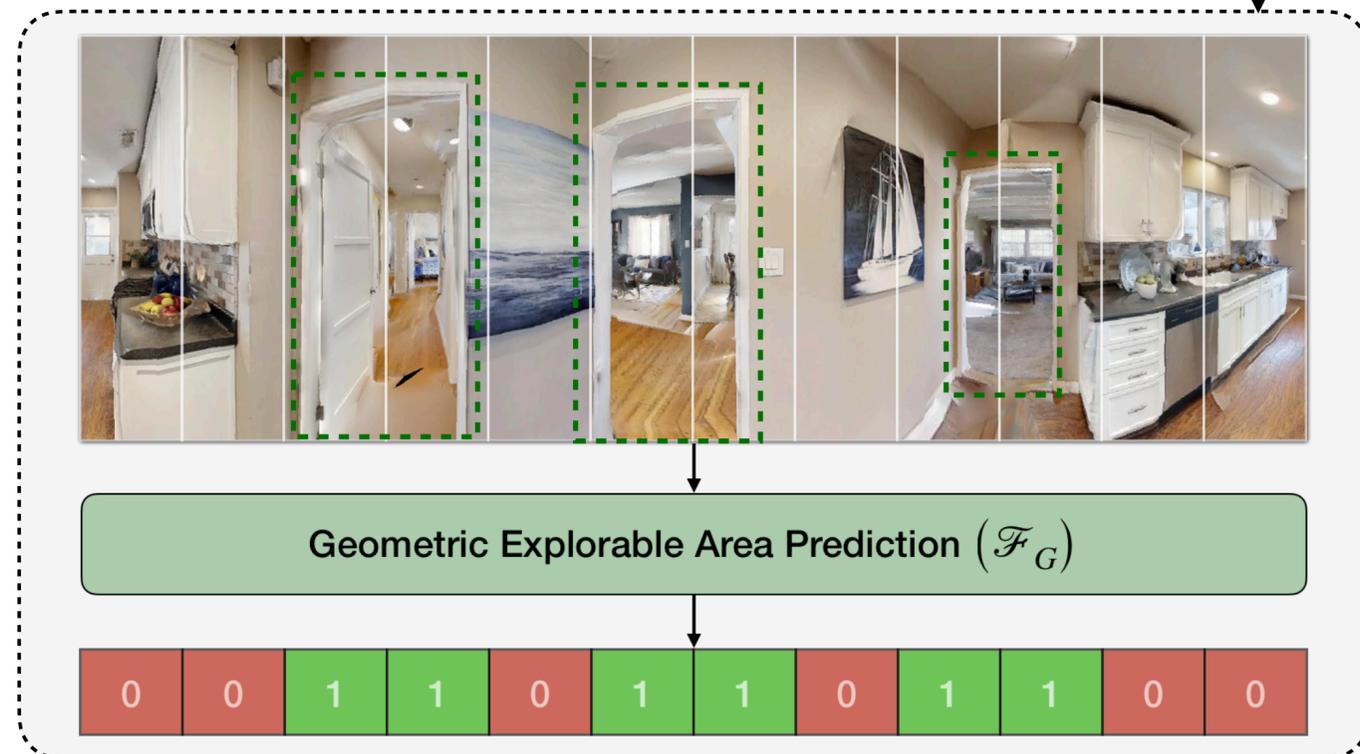
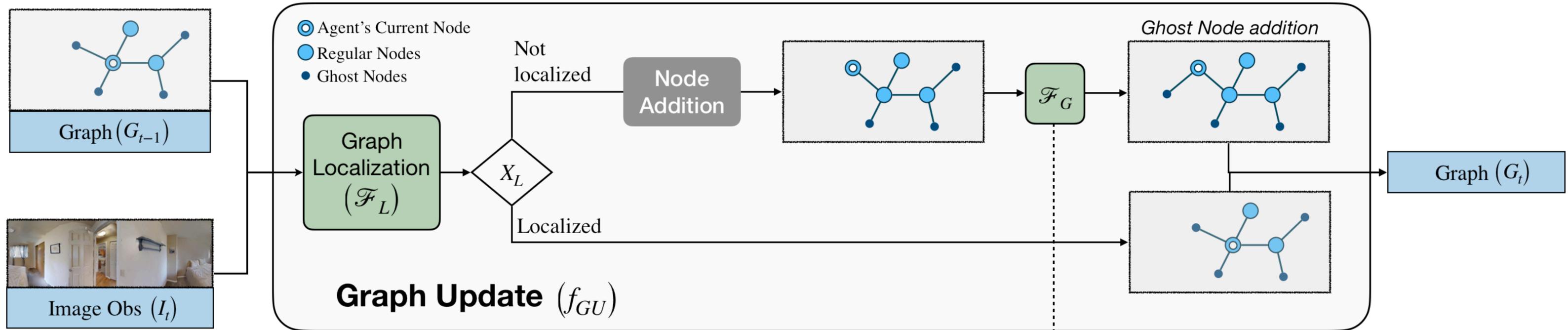
Building the Representation



Building the Representation



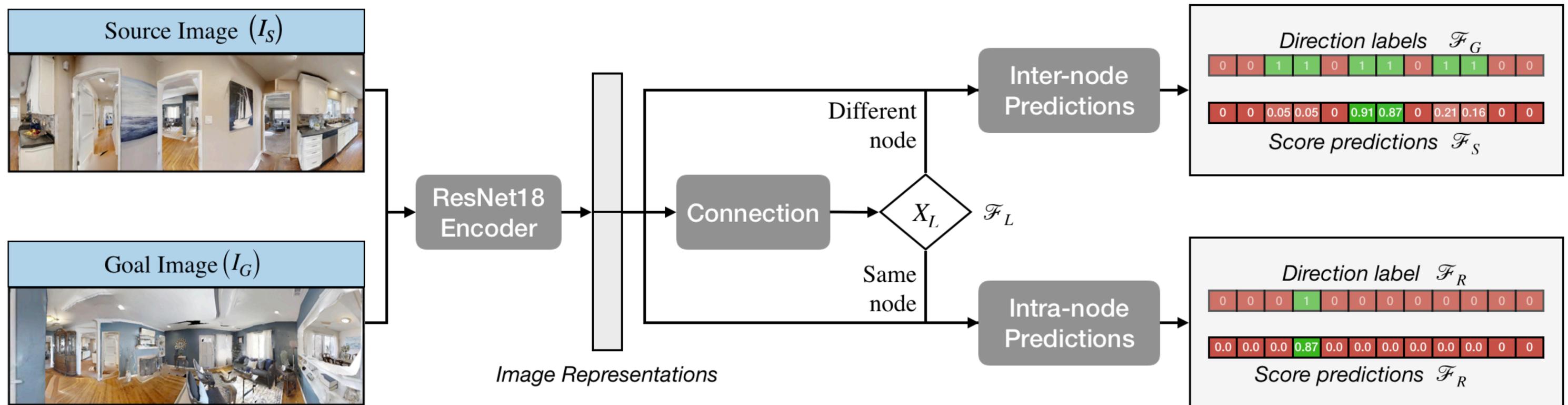
Building the Representation



Single Supervised Learning Model

- $F_G(I_1)$ = Geometric Prediction: Free directions
- $F_S(I_1, I_2)$ = Semantic Prediction: Closeness to target

- $F_R(I_1, I_2)$ = Relative Pose
- $F_L(I_1, I_2)$ = Localization



- No reinforcement learning, no interaction needed
- Can be trained completely with static data

0	0	0.05	0.05	0	0.91	0.87	0	0.21	0.16	0	0
---	---	------	------	---	------	------	---	------	------	---	---

Source Image



Target Image



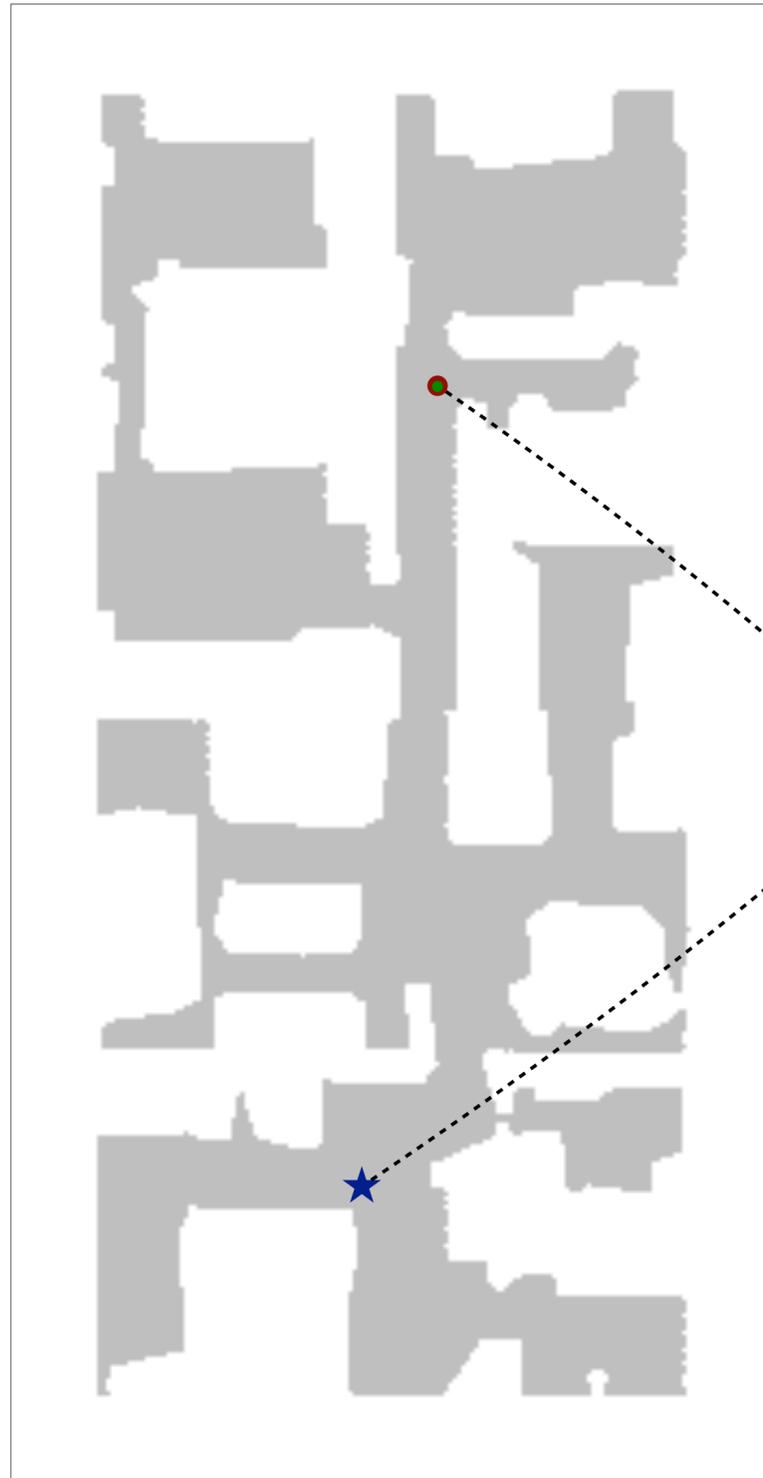
0	0	0.68	0.69	0	0.21	0.23	0	0.16	0.16	0	0
---	---	------	------	---	------	------	---	------	------	---	---

Source Image



Target Image

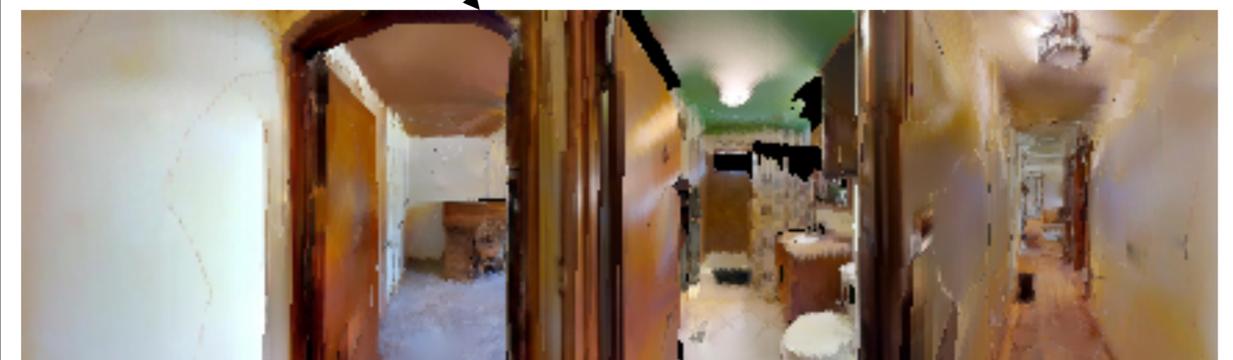




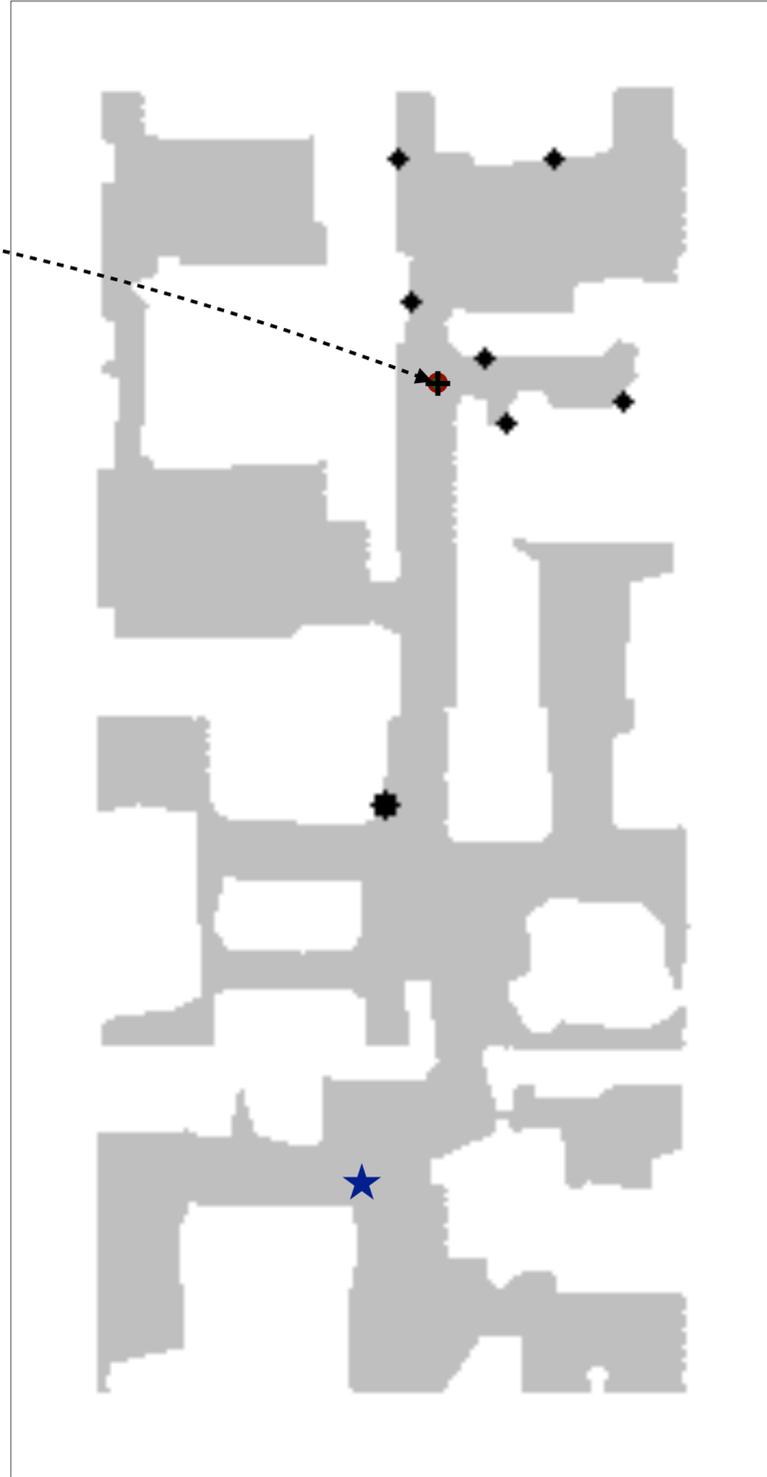
Goal Image

- Start Location** ●
- Goal Location** ★
- Current Location** ●

t = 0



Current Observation



Goal Image

- Start Location** ●
- Goal Location** ★
- Current Location** ●
- Regular Nodes** +
- Ghost Nodes** ◆
- Selected Ghost Node** ★

t = 1

Current Observation





Goal Image

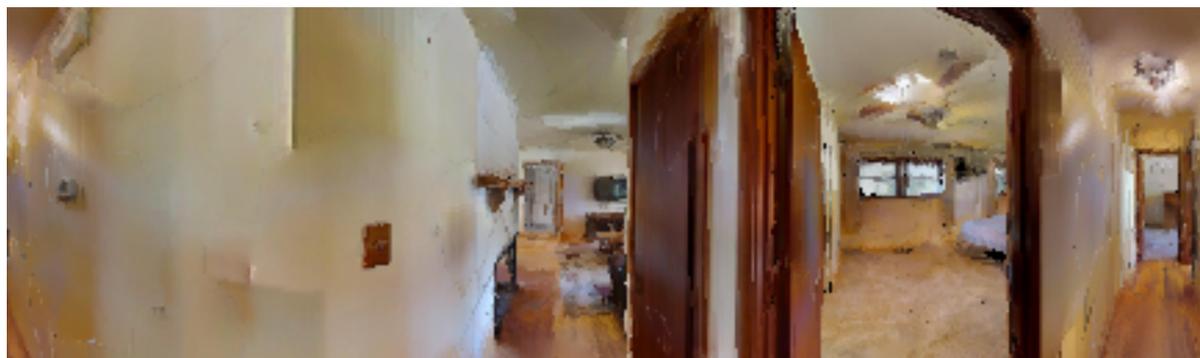


- Start Location ●
- Goal Location ★
- Current Location ●
- Regular Nodes +
- Ghost Nodes ◆
- Selected Ghost Node ■

$t = 20$

Current Observation





Goal Image

- Start Location ◆
- Goal Location ★
- Current Location ●
- Regular Nodes +
- Ghost Nodes ◆
- Selected Ghost Node ◆

t = 27

Current Observation





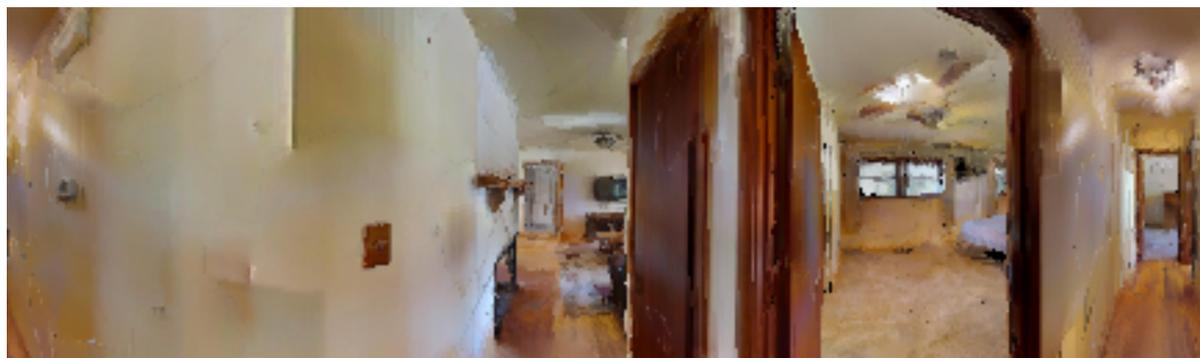
Goal Image

- Start Location ●
- Goal Location ★
- Current Location ●
- Regular Nodes +
- Ghost Nodes ◆
- Selected Ghost Node ◆

t = 27

Current Observation





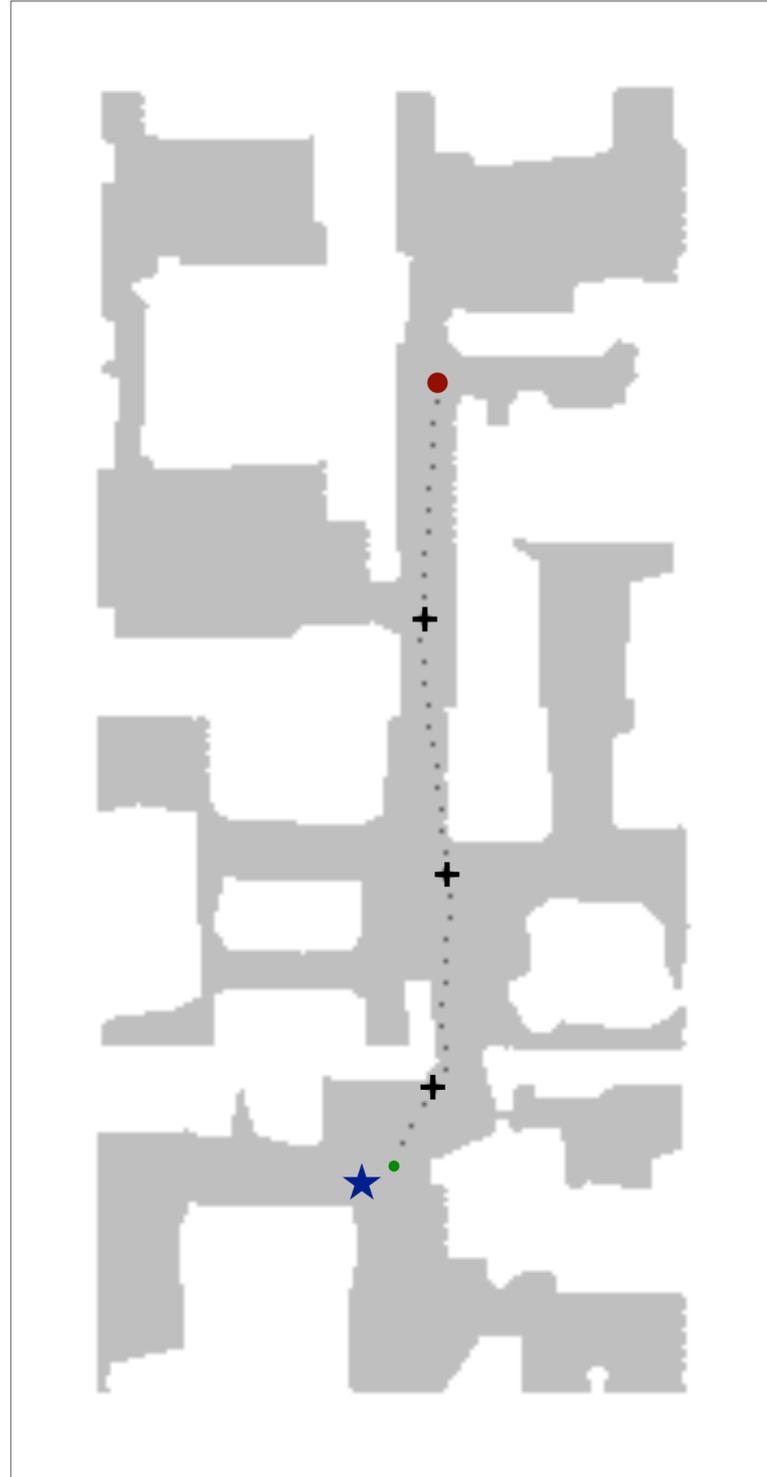
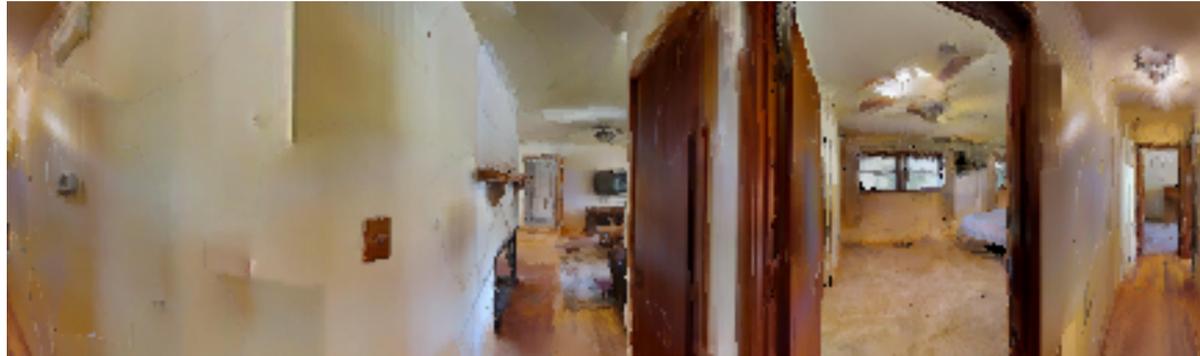
Goal Image

- Start Location ●
- Goal Location ★
- Current Location ●
- Regular Nodes +
- Ghost Nodes ◆
- Selected Ghost Node ★

t = 56

Current Observation





Goal Image

- Start Location ●
- Goal Location ★
- Current Location ●
- Regular Nodes +
- Ghost Nodes ◆
- Selected Ghost Node ◆

t = 61

Current Observation



Results (SPL)

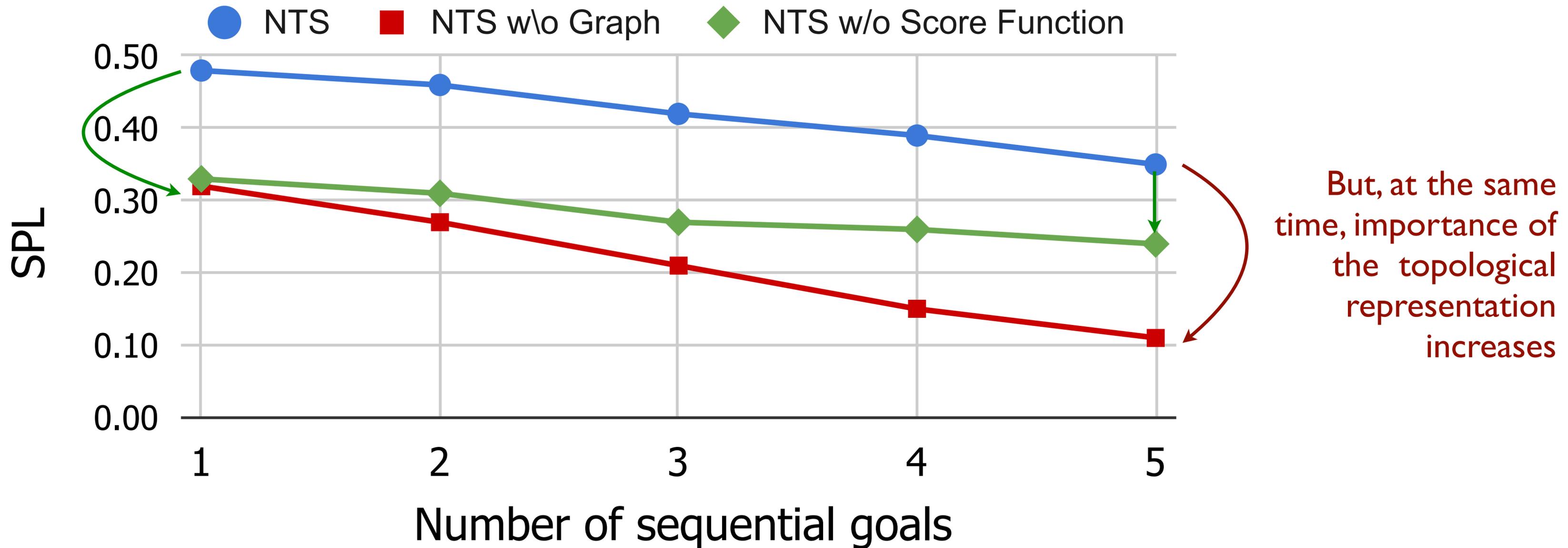
		RGB	RGBD	RGBD (No Noise)	RGBD (No Stop)	
Vanilla LSTM Memory	LSTM + Imitation	0.10	0.14	0.15	0.18	Map based methods are better than vanilla learning methods even in presence of noise
	LSTM + RL	0.10	0.13	0.14	0.17	
Metric Maps	Occupancy Maps + FBE + RL		0.26	0.31	0.24	NTS is better than occupancy map models, captures and uses semantic priors.
	ANS	0.23	0.29	0.35	0.39	
Topological Maps	NTS (Our)	0.38	0.43	0.45	0.60	

Robustness to Noise

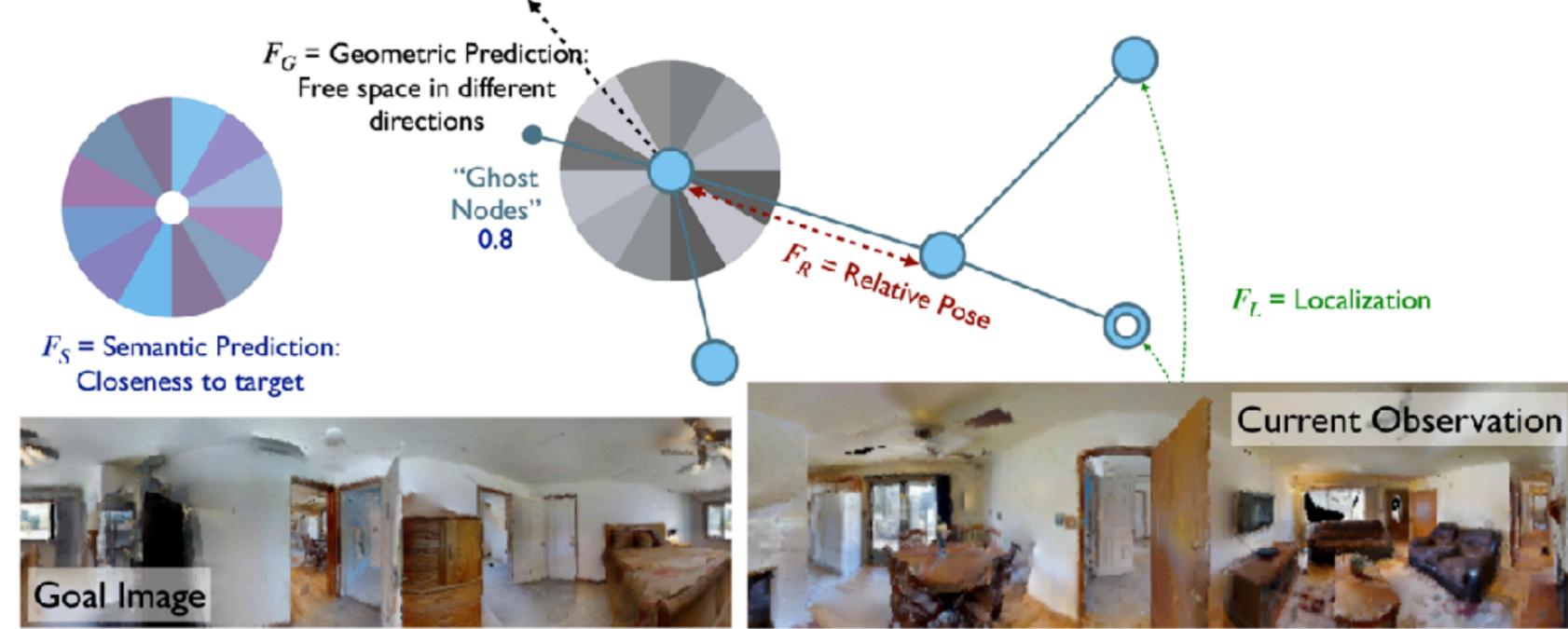
Results

Semantic score function improves efficiency when no prior experience with environment is available.

As experience in environment increases, utility of semantic function decreases



Representation for Places



- Spatial reasoning
- Semantic reasoning

- ~~Sensitive to pose error~~
- ~~Interactive training~~
- ~~Long training horizons~~

Robust to pose error

Offline supervised training

Modularized policy

- Still requires a simulator

$F_G(I_1)$: Geometry prediction

$F_R(I_1, I_2)$: Relative Pose

$F_L(I_1, I_2)$: Localization

$F_S(I_1, I_2)$: Semantic Prediction

Can we simplify and scale-up training further?

Learning F_s by Watching YouTube Videos

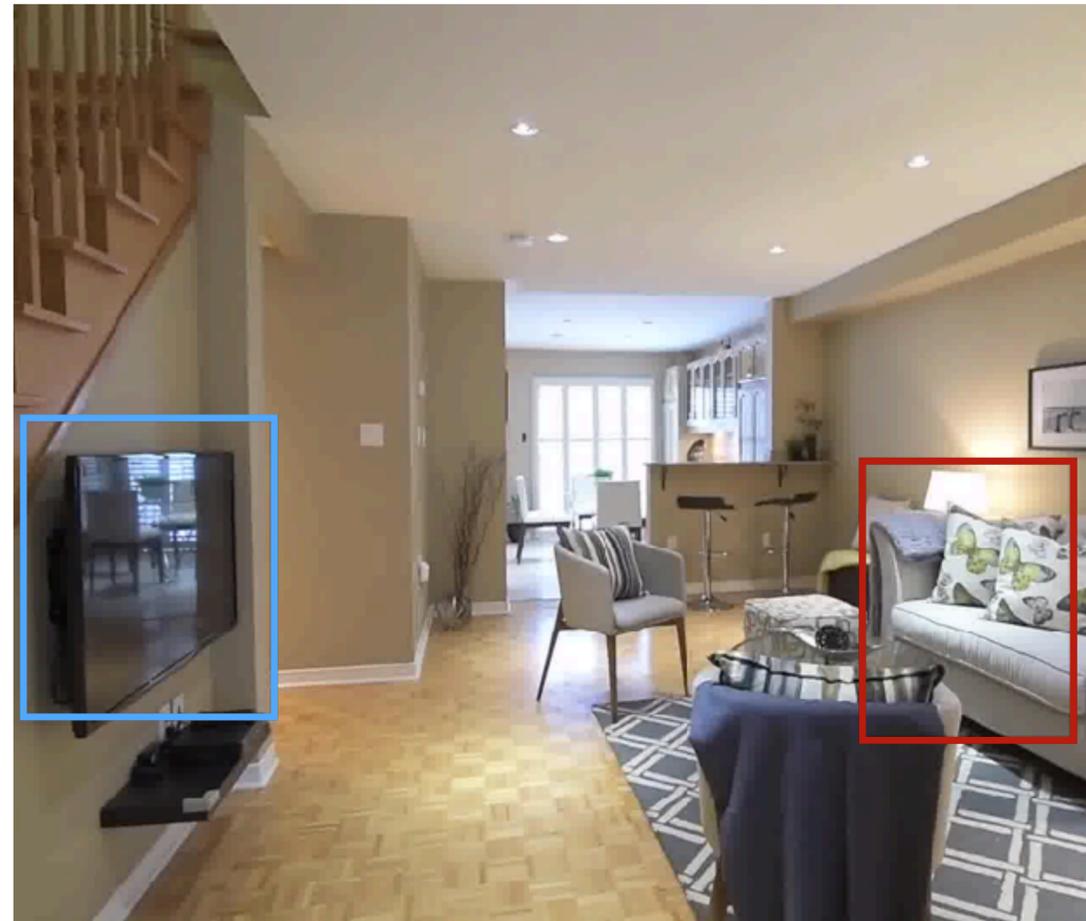
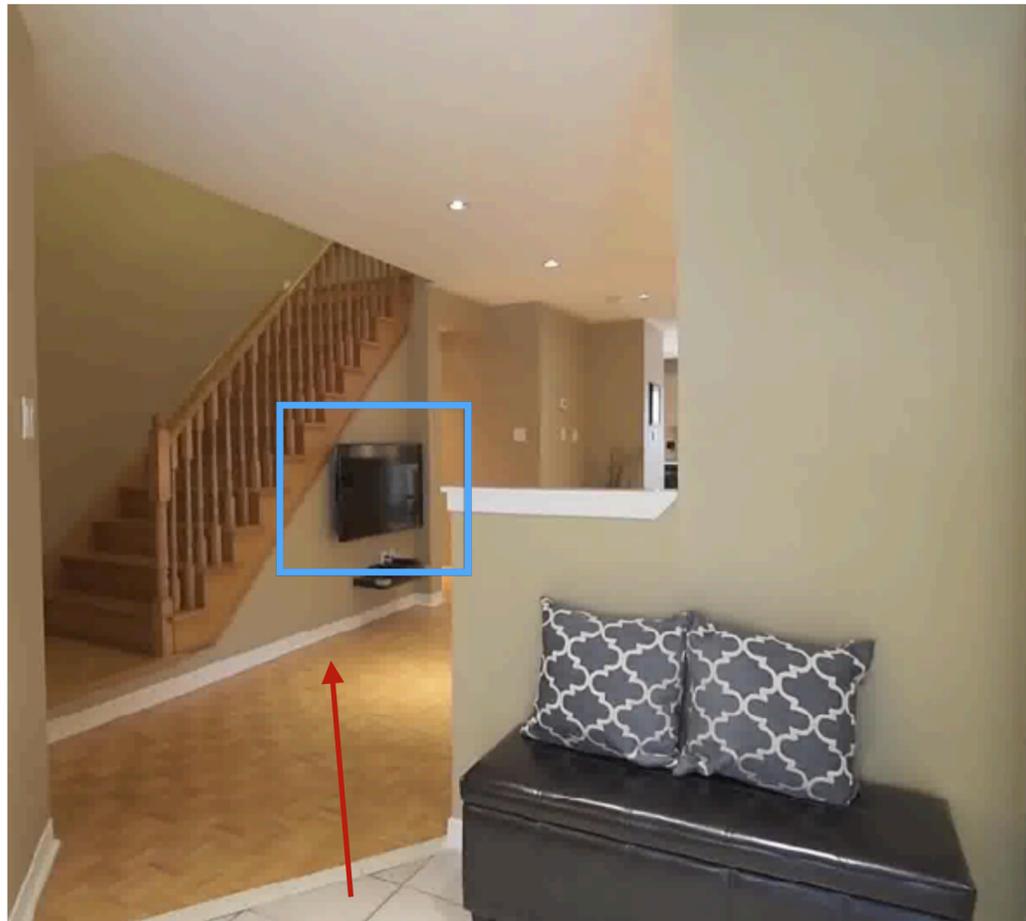


Basic Intuition

Mine for spatial co-occurrences

Video

time



e.g. cues for finding a couch

Challenges in Using Such Videos

- Videos don't come with action labels
 - ⇒ Action Grounding via an Inverse Model [1]
- Goals and intents are not known
 - ⇒ Use off-the-shelf object detectors to label frames with desired objects
- Depicted trajectories may not be optimal
 - ⇒ Use Q-learning to learn optimal behavior from sub-optimal data [2]

[1] A. Kumar, S. Gupta, J. Malik. Learning navigation subroutines by watching videos. In *CoRL*, 2019.

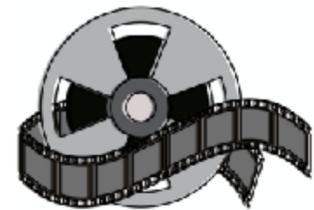
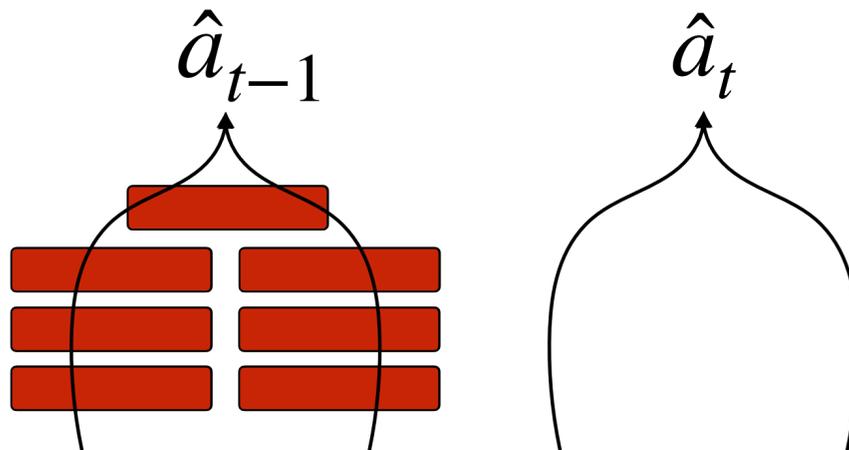
[2] Watkins, C. J. C. H. (1989). Learning from delayed rewards.

Value Learning from Videos (VLV)

a) Action Grounding

Inverse Model

built by executing random actions on robot

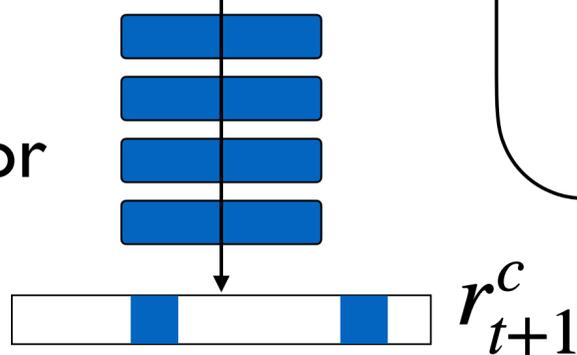


Real Estate Tour from YouTube



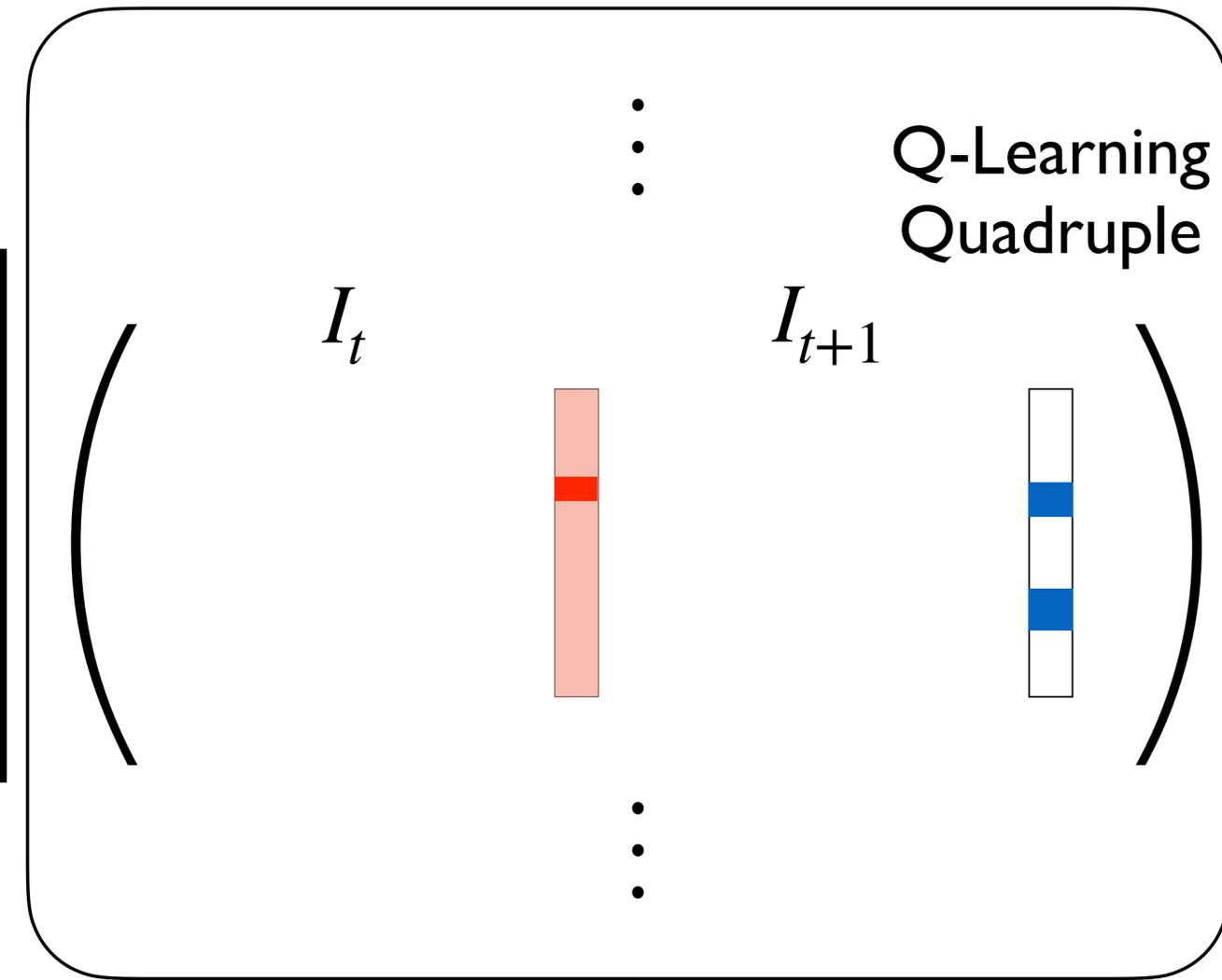
Object Detector

trained on COCO



b) Goal Labeling

Value function that uses implicitly learns semantic cues for seeking objects in novel indoor environments



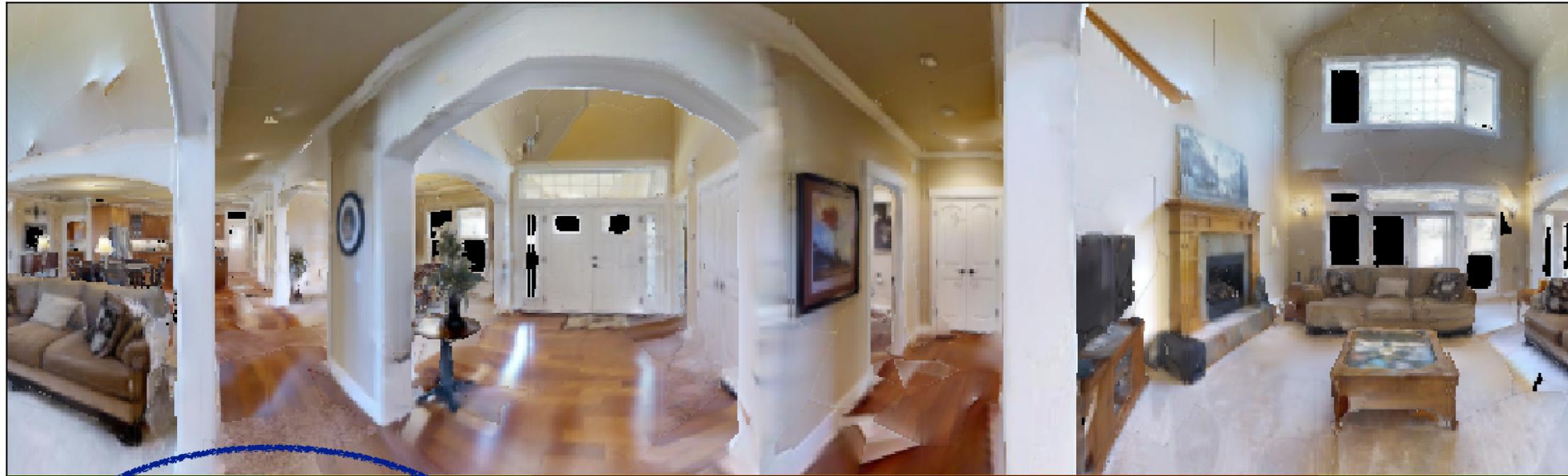
c) Q-Learning

$$\rightarrow f(I, c) = \max_a Q^*(I, a, c)$$

Learned Value Function

$$f(I, c) \approx \text{nearness to goal}$$

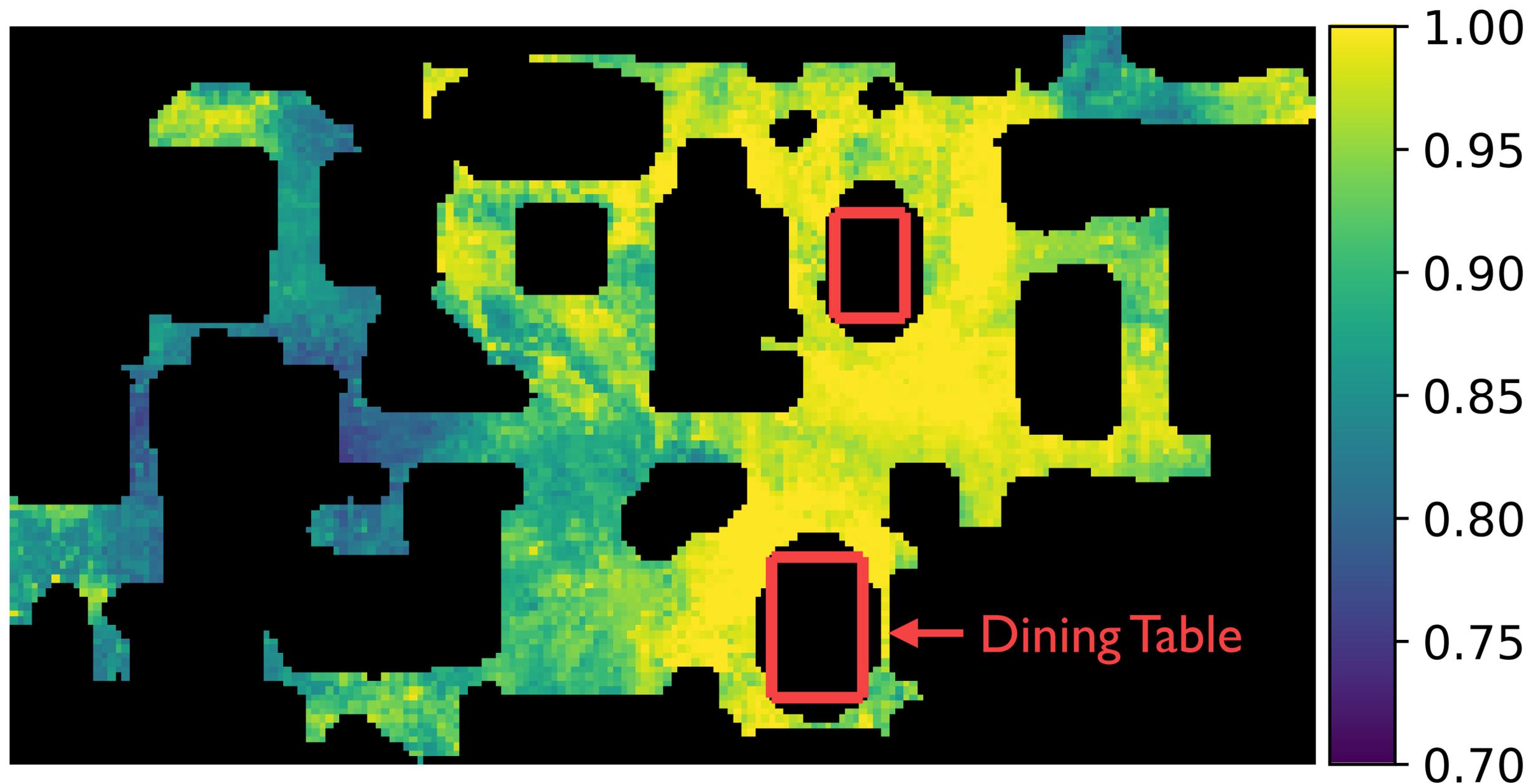
Value function predicts a proxy for nearness to a goal object for a given image



Learned Value Function

$$f(I, c) \approx \text{nearness to goal}$$

Value function predicts a proxy for nearness to a goal object for a given image

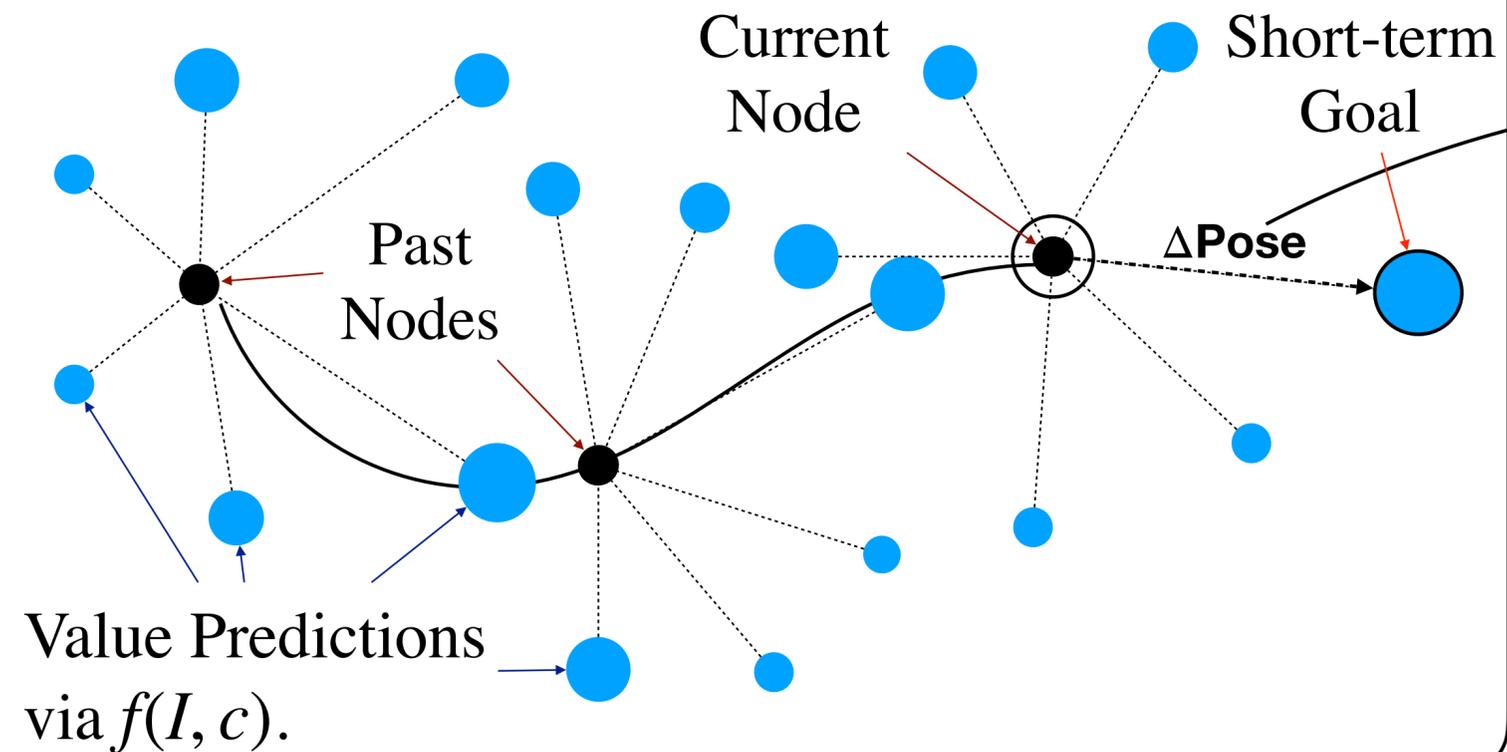


Using Learned Values for Semantic Navigation

Hierarchical Policy

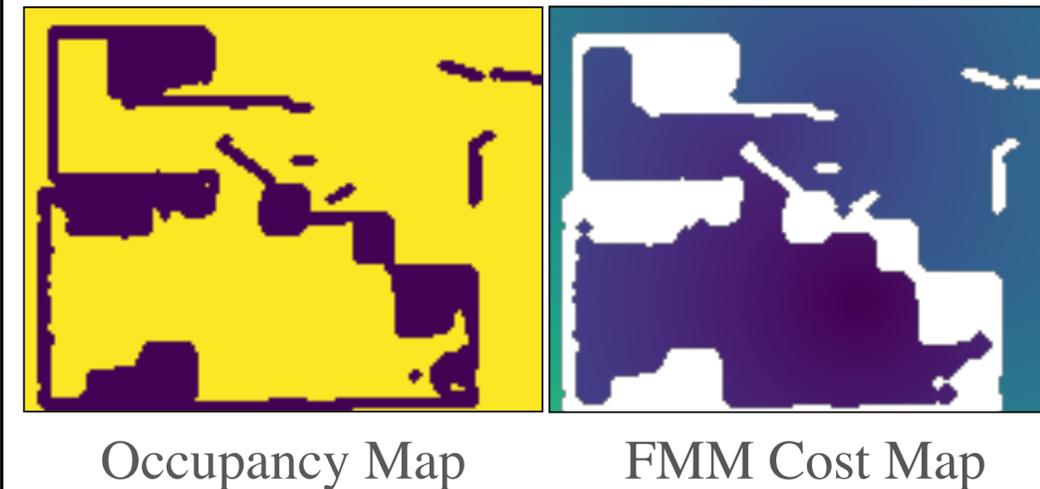
High-Level Policy

- Decides where to go next and emits short-term goal
- Builds a topological map [1] that stores values predicted by $f(I, c)$ at different locations in different directions



Low-Level Policy

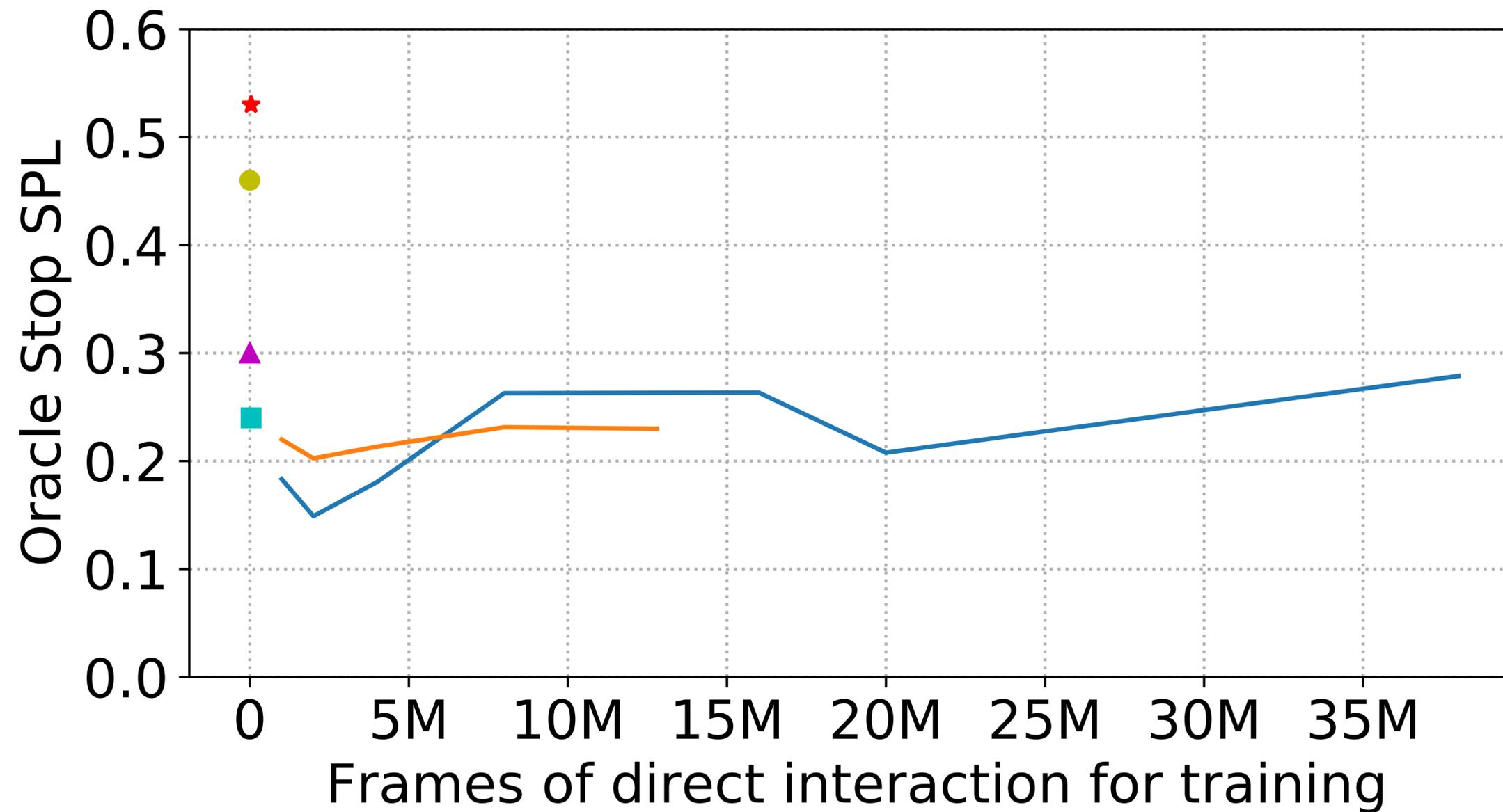
- Executes actions to achieve short-term goal
- Incrementally builds occupancy map from depth camera, plans paths



Forward
Left
Right
Stop

[1] D. Chaplot, R. Salakhutdinov, A. Gupta, S. Gupta. Neural topological slam for visual navigation. In *CVPR*, 2020.

Results



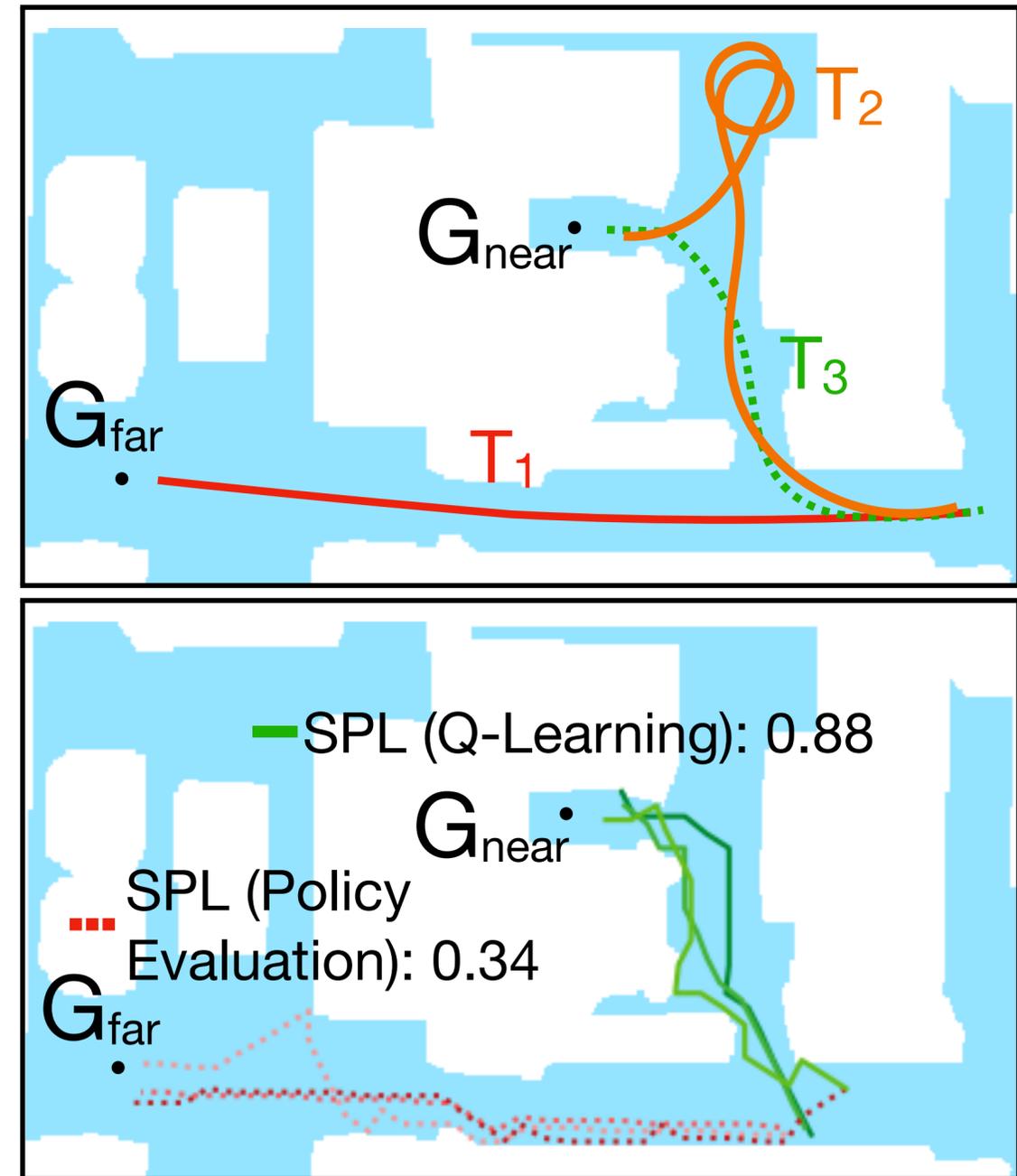
▲ [0.30] Topological Exploration
● [0.46] Detection Seeker
■ [0.24] Behavior Cloning (YouTube)

★ [0.53] Ours (YouTube)
— [0.28] RL
— [0.23] BC (YouTube) + RL

Ablations

Method	Oracle Stop SPL			
	Easy	Medium	Hard	Overall
Base Setting	0.62	0.42	0.23	0.40
True Actions	0.61	0.45	0.25	0.41
True Detections	0.62	0.45	0.22	0.40
True Rewards	0.64	0.46	0.21	0.41
Optimal Trajectories	0.65	0.46	0.25	0.43
Detector Score	0.73	0.48	0.26	0.46
Train on 360° Videos	0.66	0.51	0.32	0.47
No Hierarchy	0.38	0.10	0.02	0.15

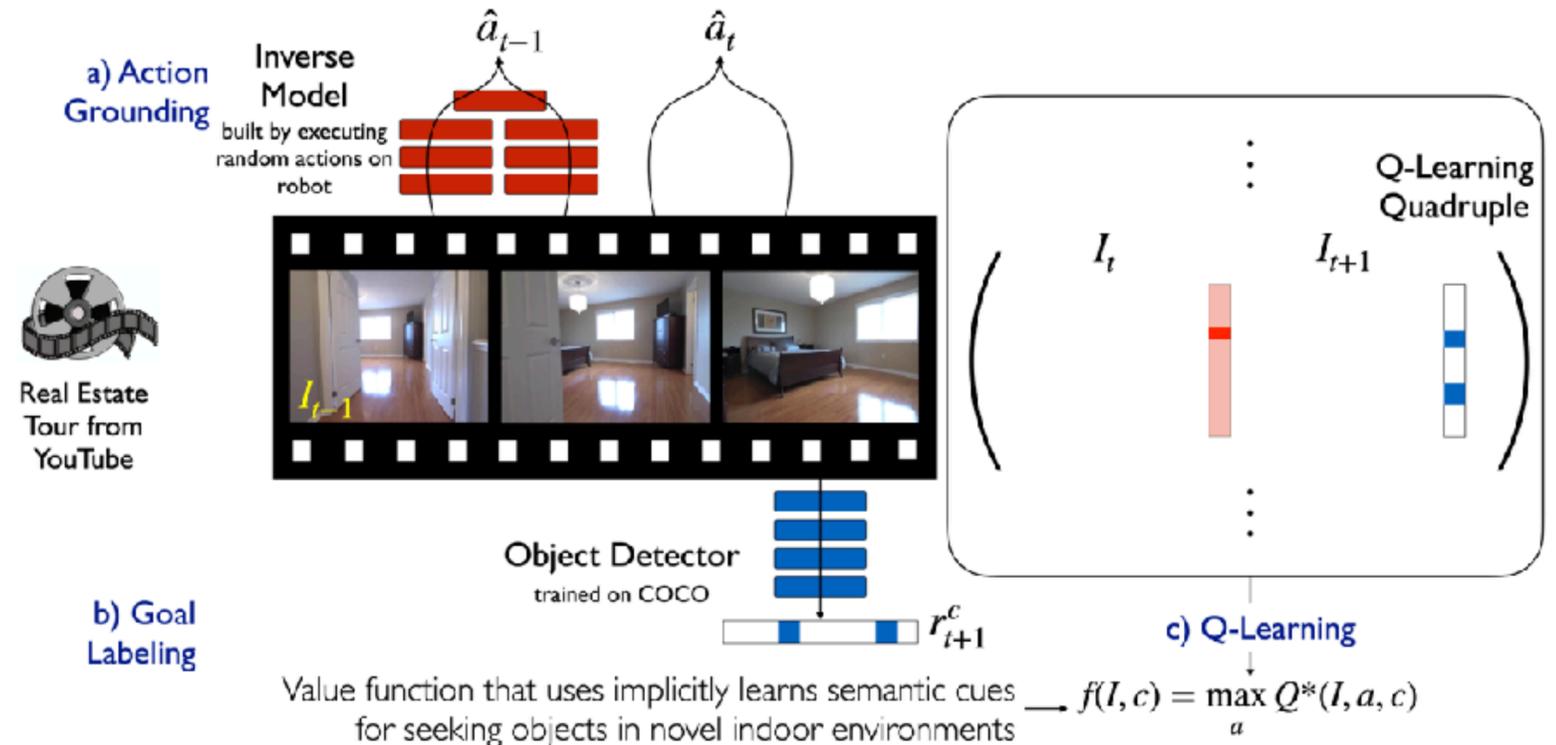
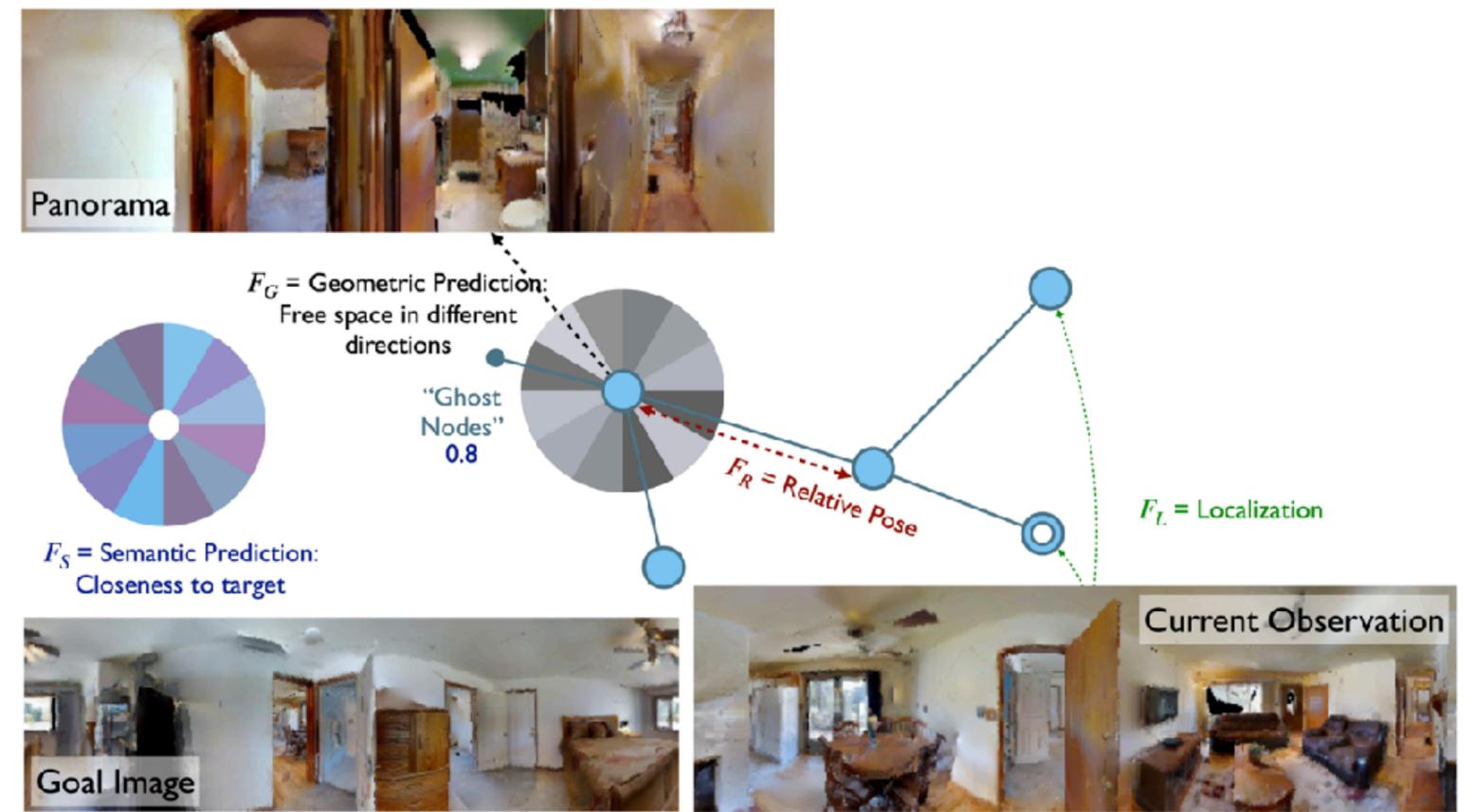
- Inverse model and detector do not hurt performance significantly
- Detector at test time helps for close objects, panorama helps for far objects
- Q-Learning outperforms simple policy evaluation for challenging environments
- Hierarchical policy is a major factor in strong performance



Navigation to couches in novel environments

Representation for Places

- Spatial reasoning
- Semantic reasoning
- Robust to pose error
- Modularized policy
- Training from in-the-wild videos



In this talk,

Representations for Places that Afford Navigation in Novel Environments

- *Augmenting metric representations with semantic reasoning*
- *Relaxing the need for metric representations*
- *Scaling-up training of such representations*

Operationalize insights from classical robotics into learning paradigms