

SEAL: Self-supervised Embodied Active Learning

NeurIPS 2021

Webpage: <https://devendrachaplot.github.io/projects/seal>



**Devendra Singh
Chaplot**



**Murtaza
Dalal**



**Saurabh
Gupta**



**Jitendra
Malik**



**Ruslan
Salakhutdinov**

Internet Computer Vision

Internet Data



[1]

Semantic Segmentation



GRASS, CAT,
TREE, SKY

Classification + Localization



CAT

Object Detection



DOG, DOG, CAT

Instance Segmentation



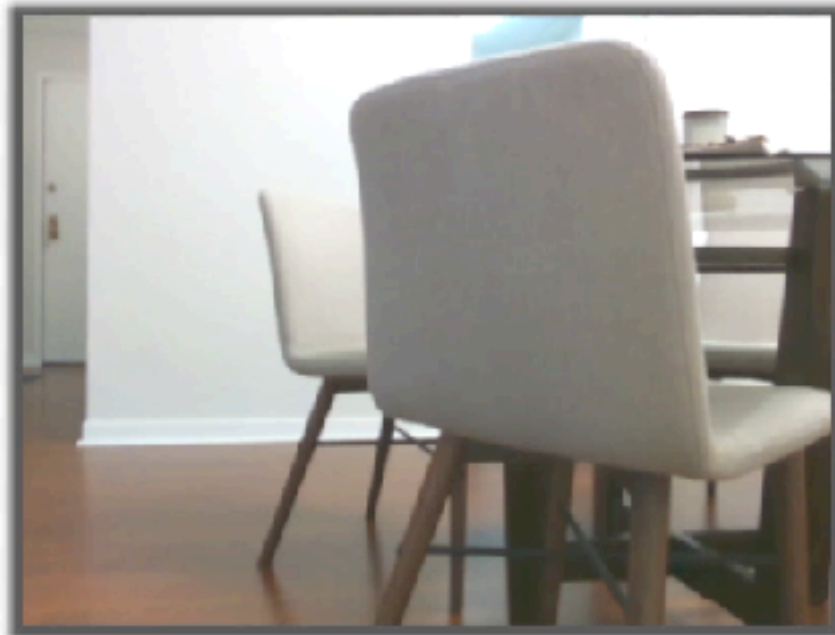
DOG, DOG, CAT [2]

[1] Karpathy. <https://cs.stanford.edu/people/karpathy/cnnembed/>

[2] Li, Johnson, Yeung. http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture11.pdf

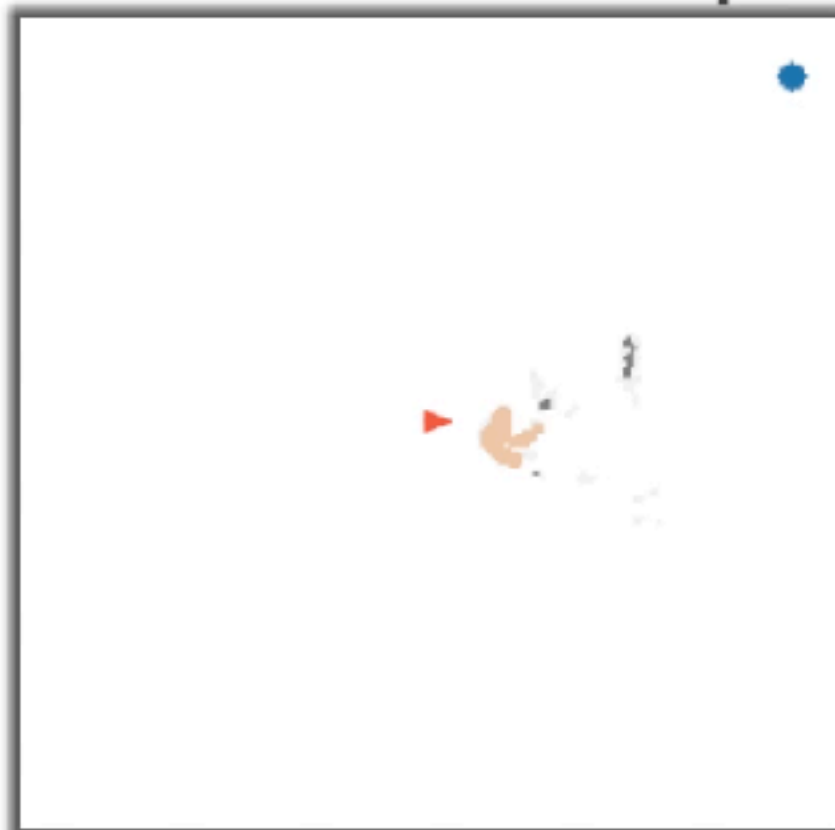
Embodied Agents

Observation

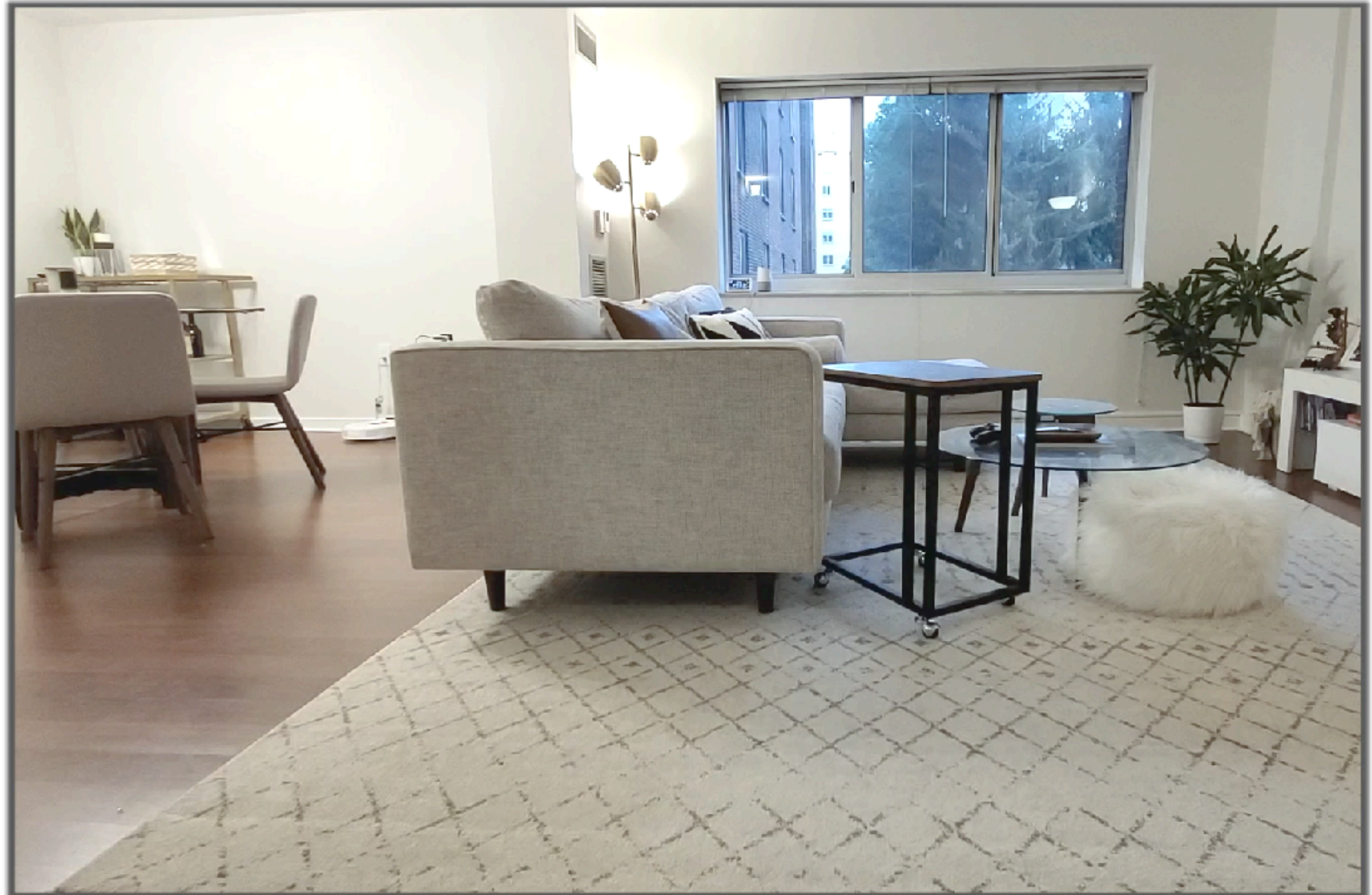


Goal: *Potted Plant*

Predicted
Semantic Map



Third-person view



[Chaplot et al. Object Goal Navigation using Goal-Oriented Semantic Exploration. NeurIPS-20]

ROBOT VISION



Berthold Klaus Paul Horn

Contents

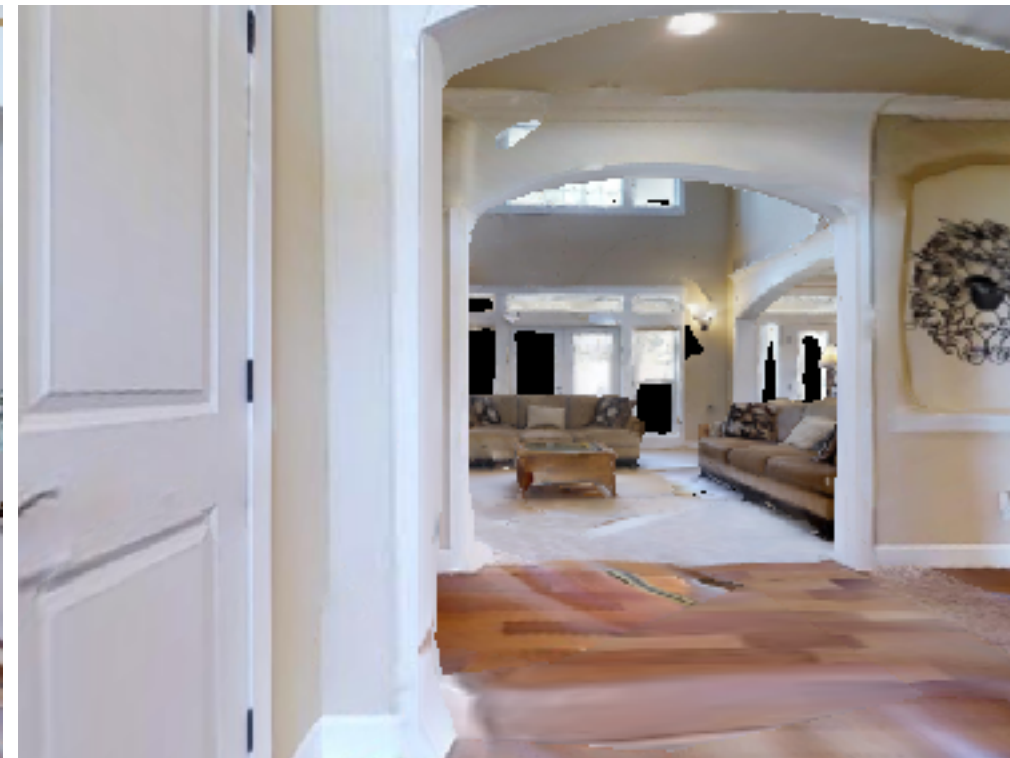
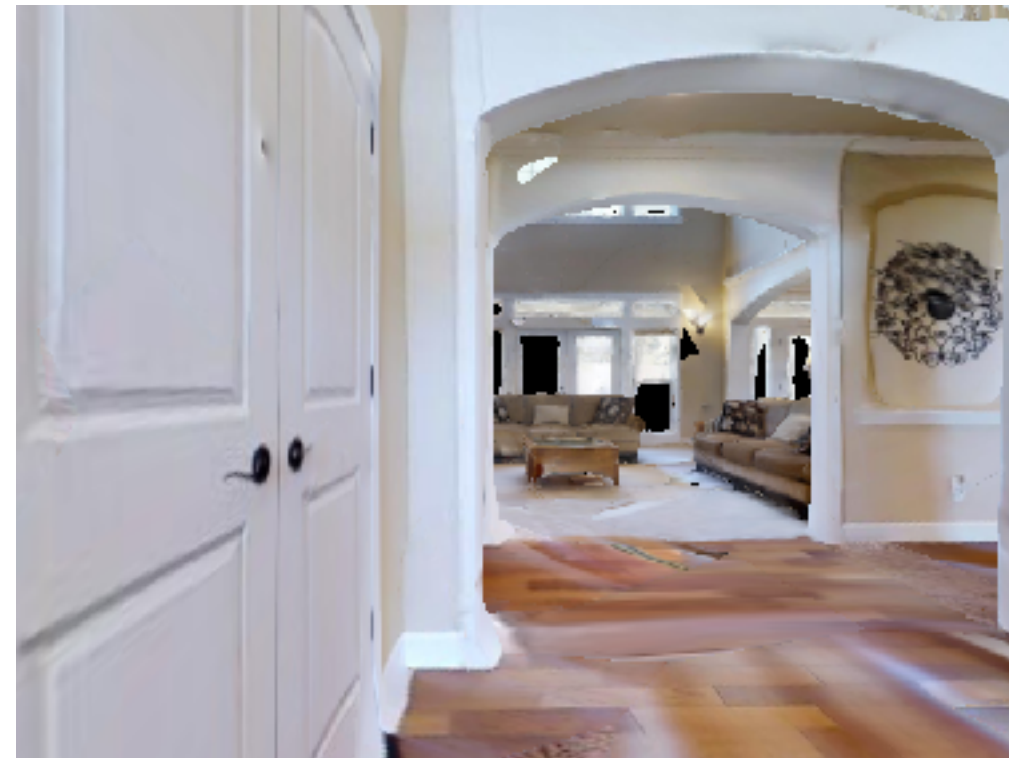
Preface	vii
Acknowledgments	xi
1 Introduction	1
2 Image Formation & Image Sensing	18
3 Binary Images: Geometrical Properties	46
4 Binary Images: Topological Properties	65
5 Regions & Image Segmentation	90
6 Image Processing: Continuous Images	103
7 Image Processing: Discrete Images	144
8 Edges & Edge Finding	161
9 Lightness & Color	185
10 Reflectance Map: Photometric Stereo	202
11 Reflectance Map: Shape from Shading	243
12 Motion Field & Optical Flow	278
13 Photogrammetry & Stereo	299
14 Pattern Classification	334
15 Polyhedral Objects	349
16 Extended Gaussian Images	365
17 Passive Navigation & Structure from Motion	400
18 Picking Parts out of a Bin	423
Appendix: Useful Mathematical Techniques	453
Bibliography	475
Index	503

Internet vs Embodied Data

Static Internet data



Active Embodied data



Using Internet models for Embodied Agents



False positives



False negatives

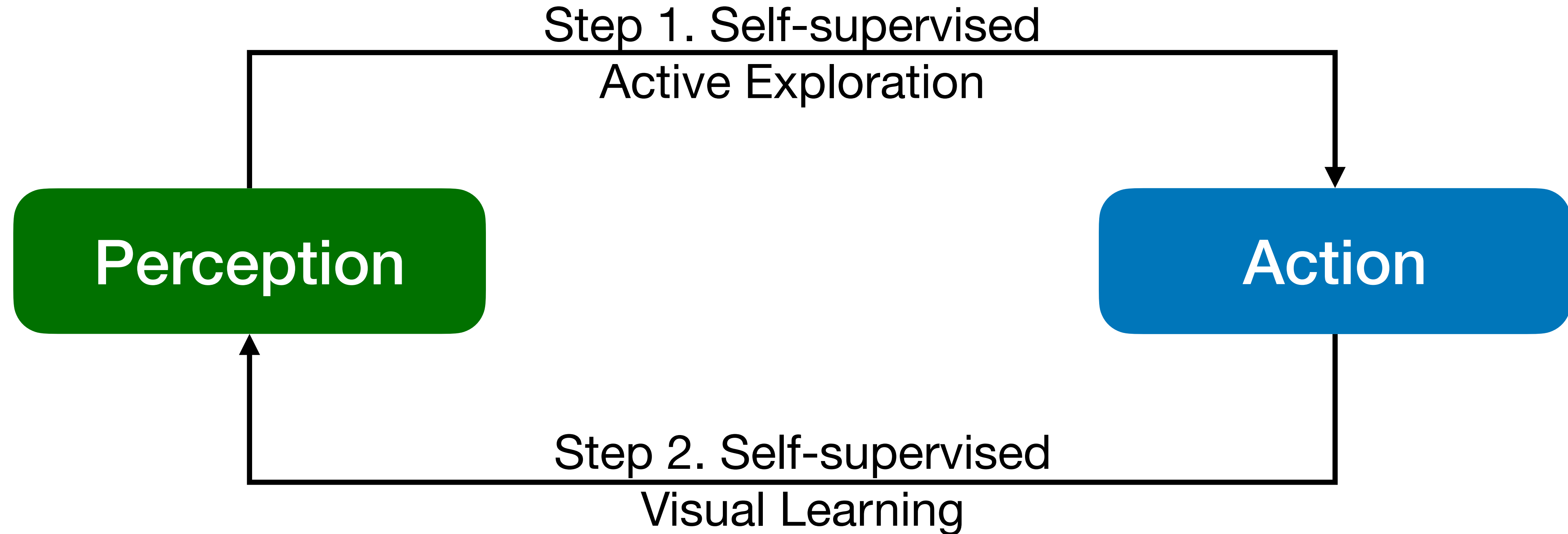
[Chaplot et al. Object Goal Navigation using Goal-Oriented Semantic Exploration. NeurIPS-20]

Embodied Perception

Active Embodied data

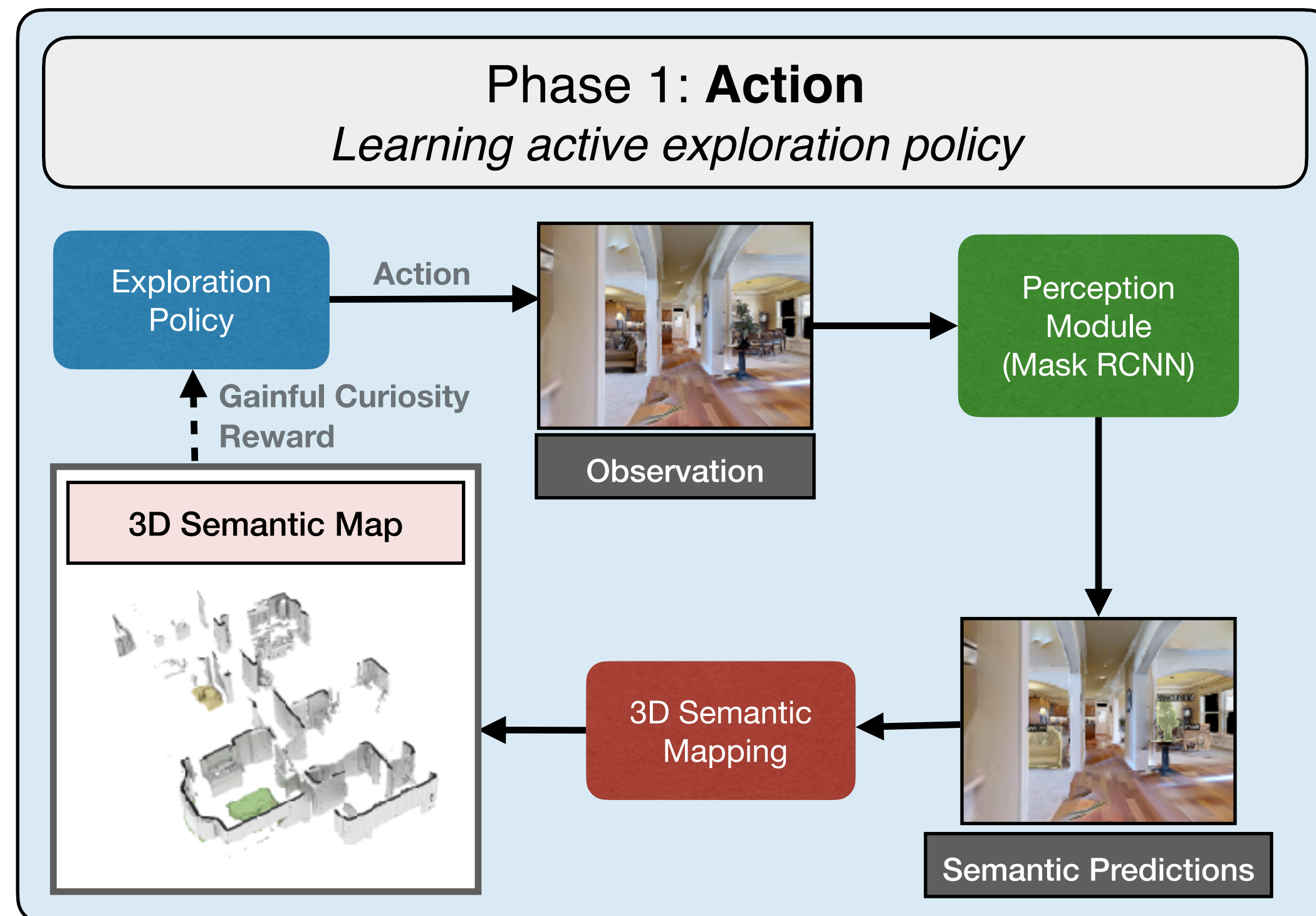


Perception-Action Loop

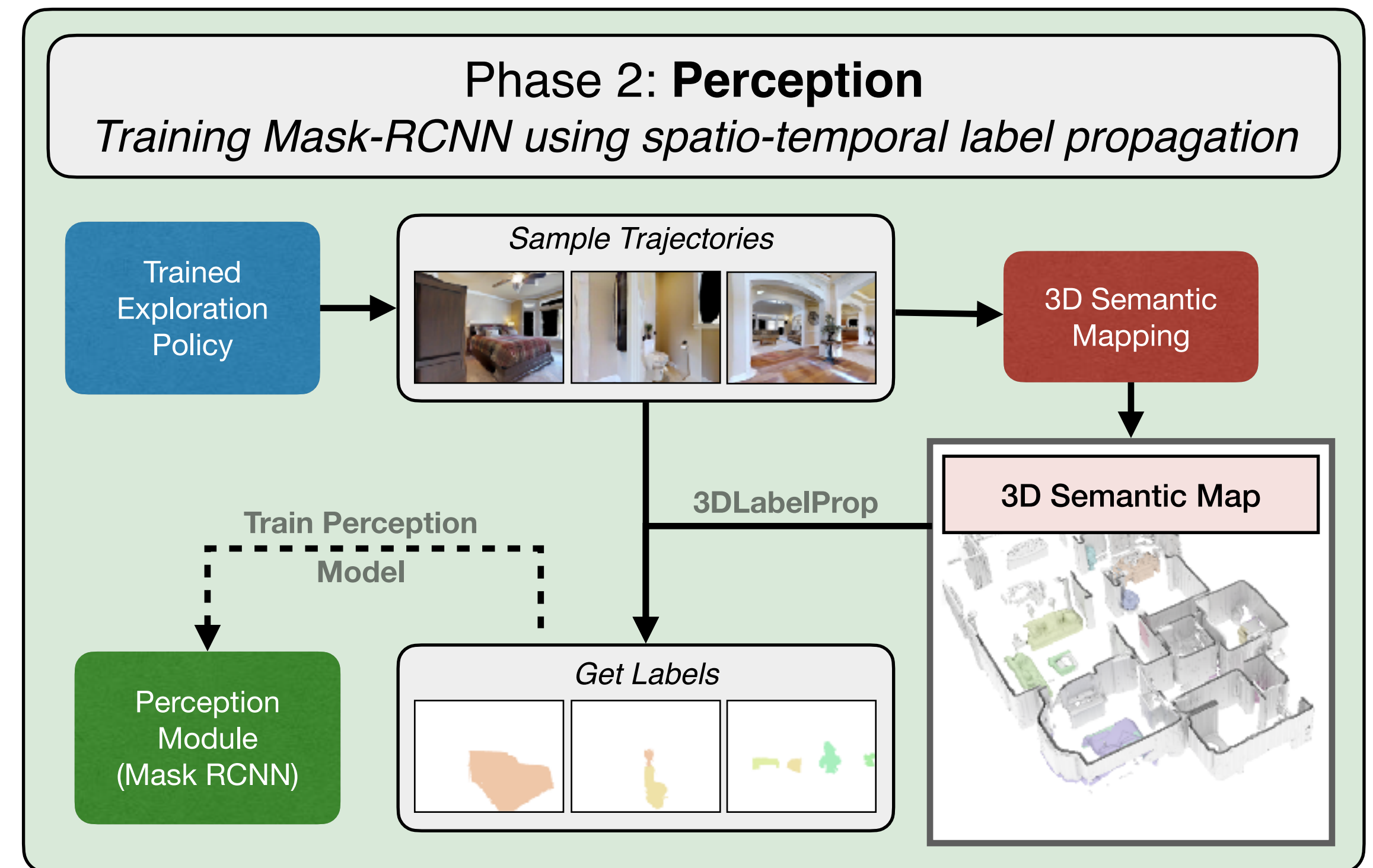
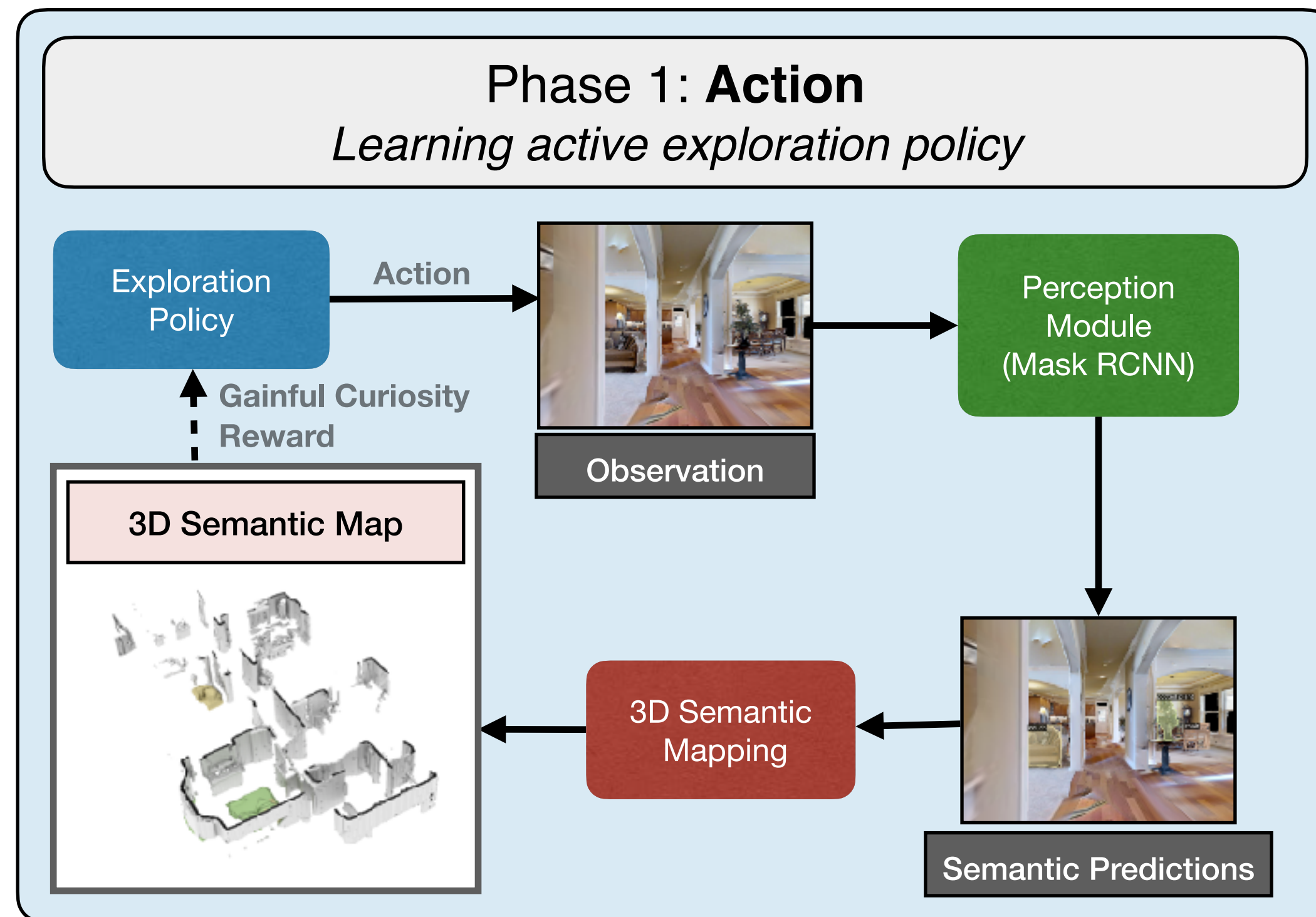


We must perceive in order to move, but we must also move in order to perceive
- Gibson (1979)

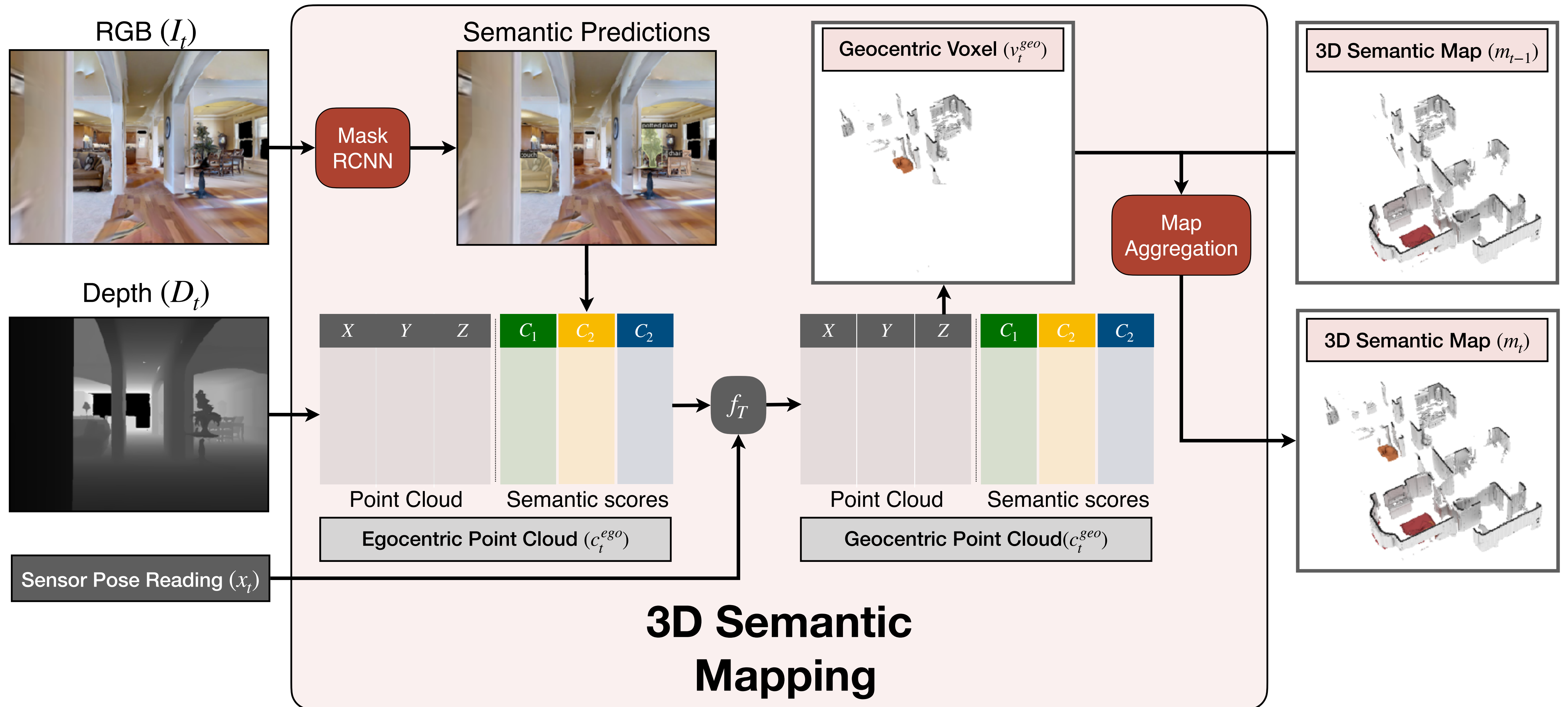
SEAL: Self-supervised Embodied Active Learning



SEAL: Self-supervised Embodied Active Learning



3D Semantic Mapping



3D Semantic Mapping

RGB Observation



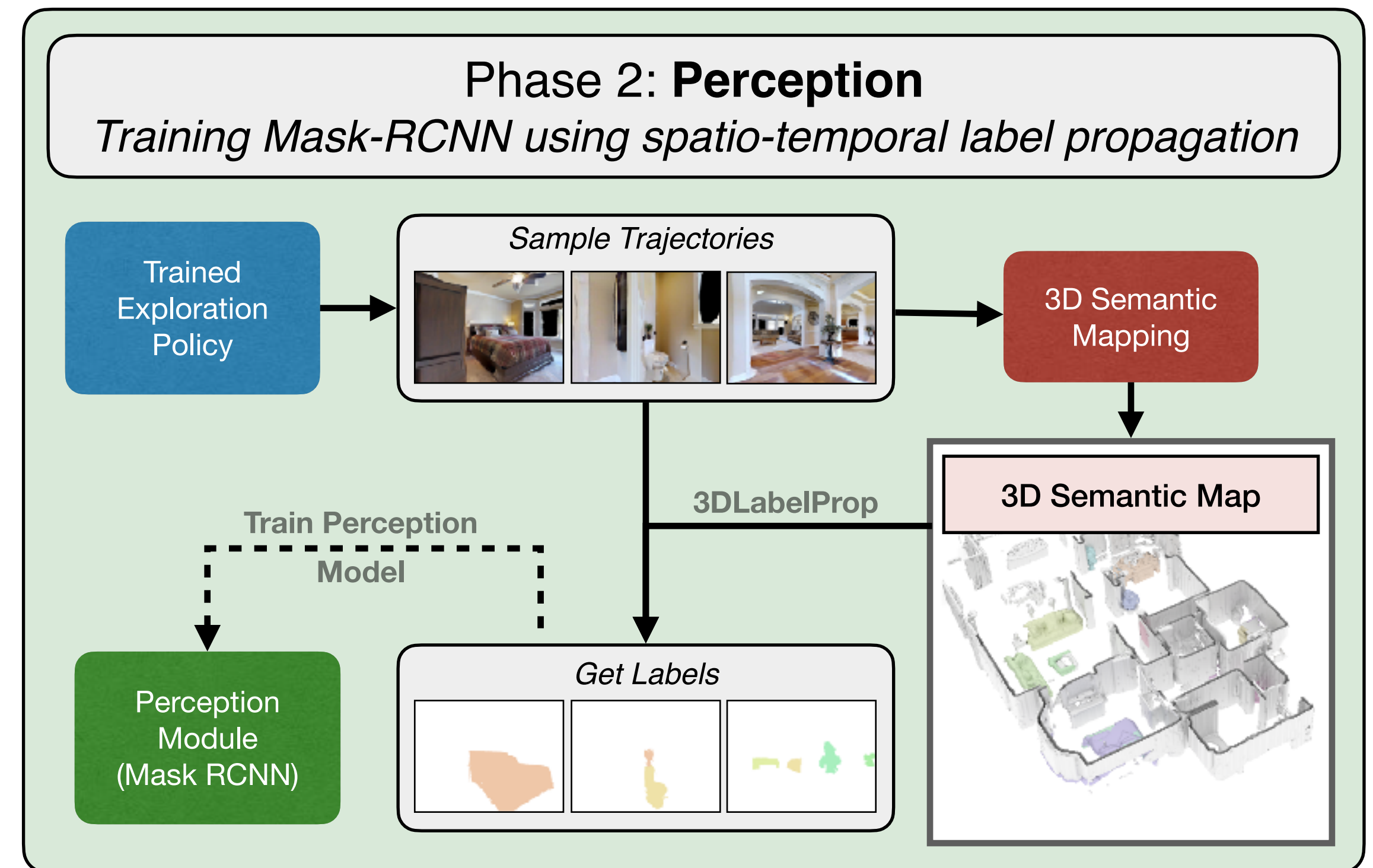
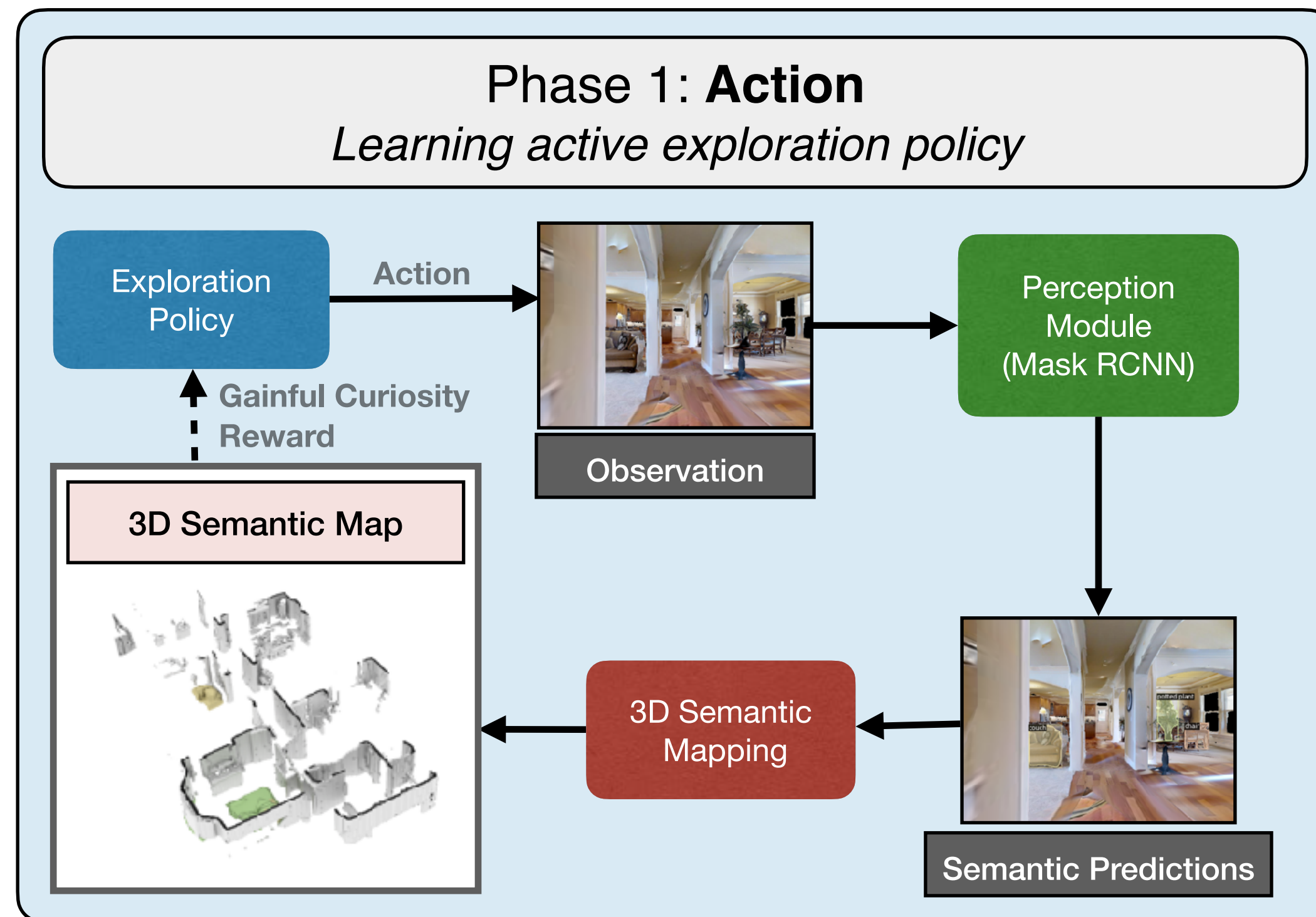
3D Semantic Map



Mask-RCNN Predictions



SEAL: Self-supervised Embodied Active Learning



3D Semantic Mapping

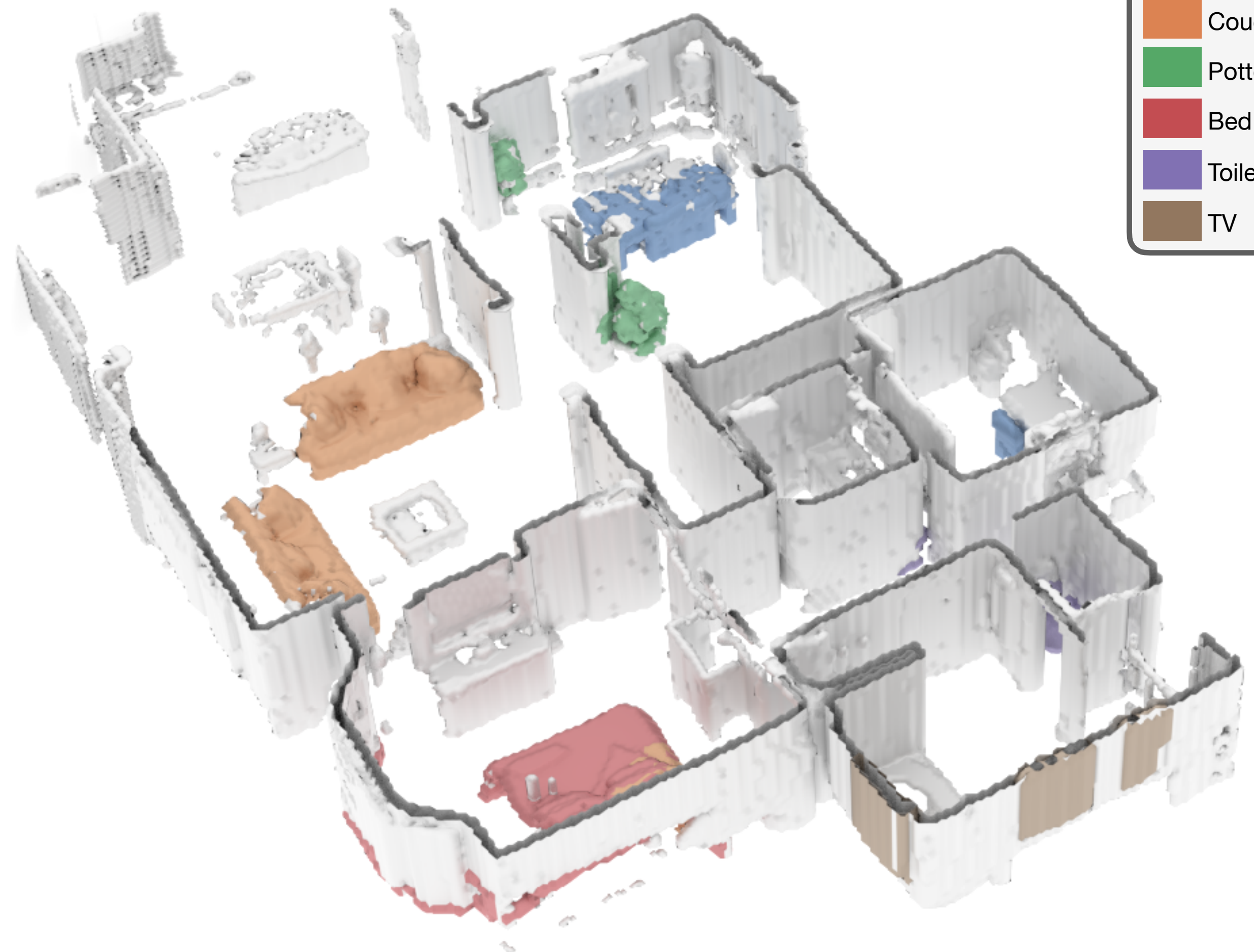
RGB Observation



Mask-RCNN Predictions

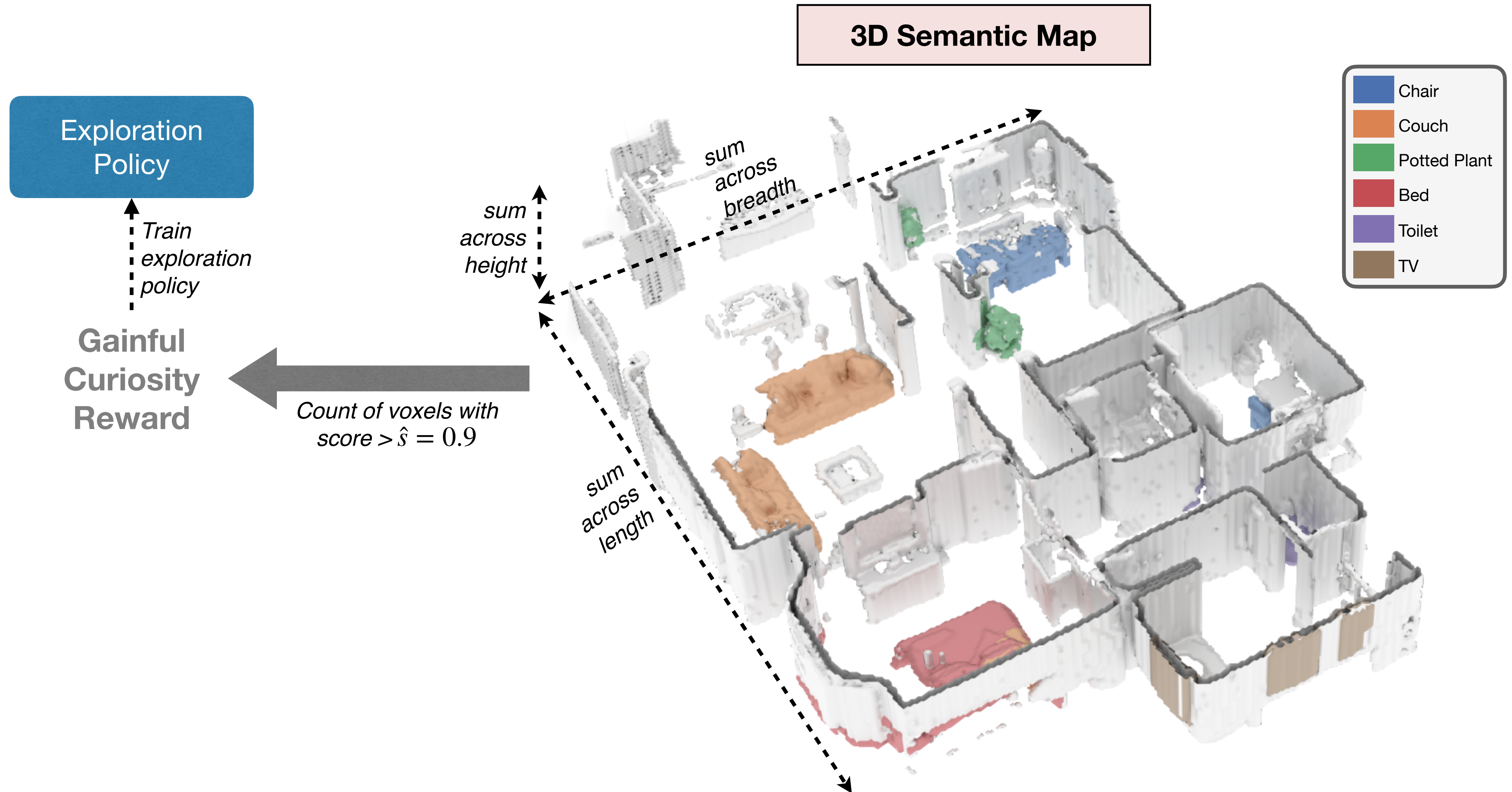


3D Semantic Map

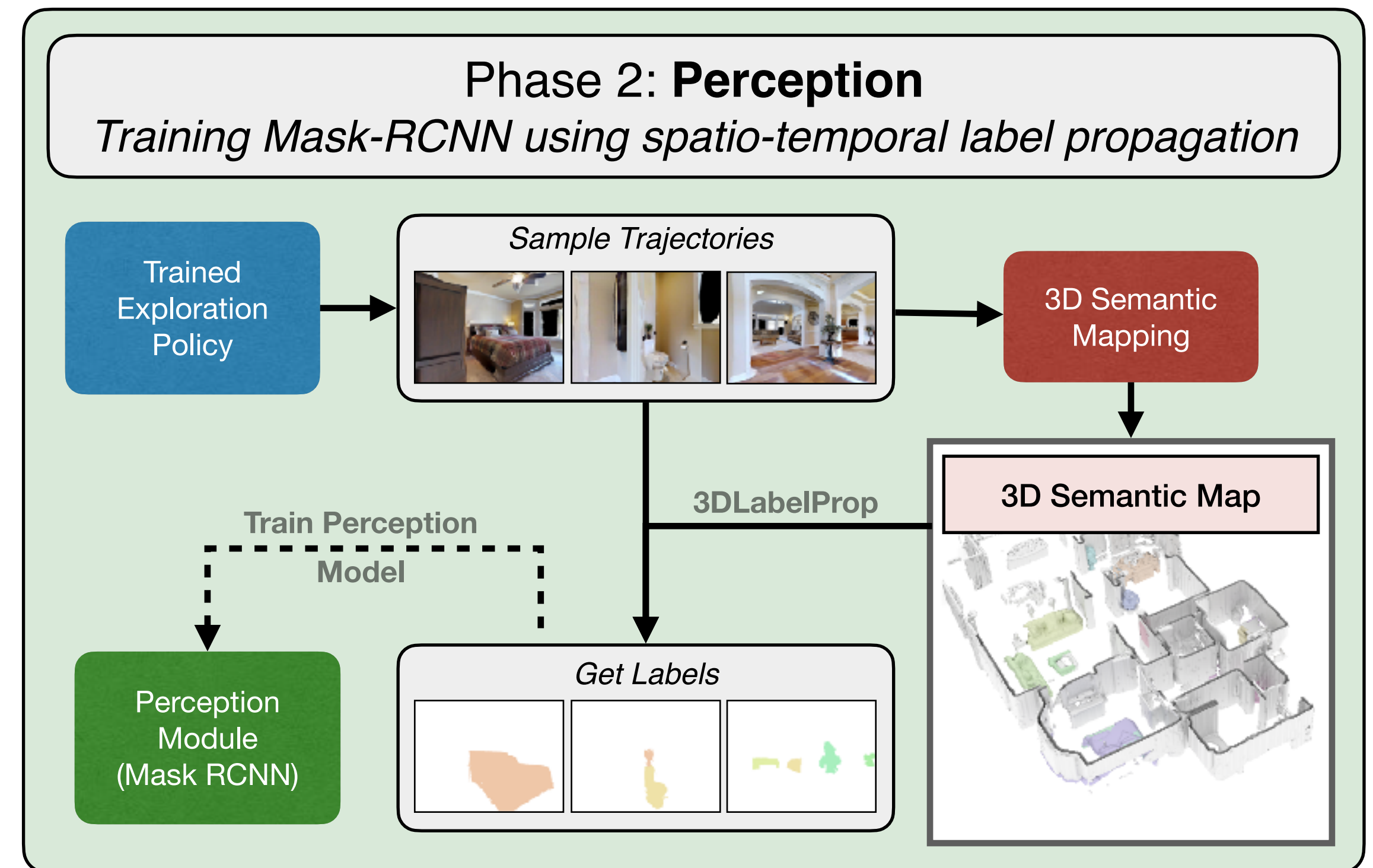
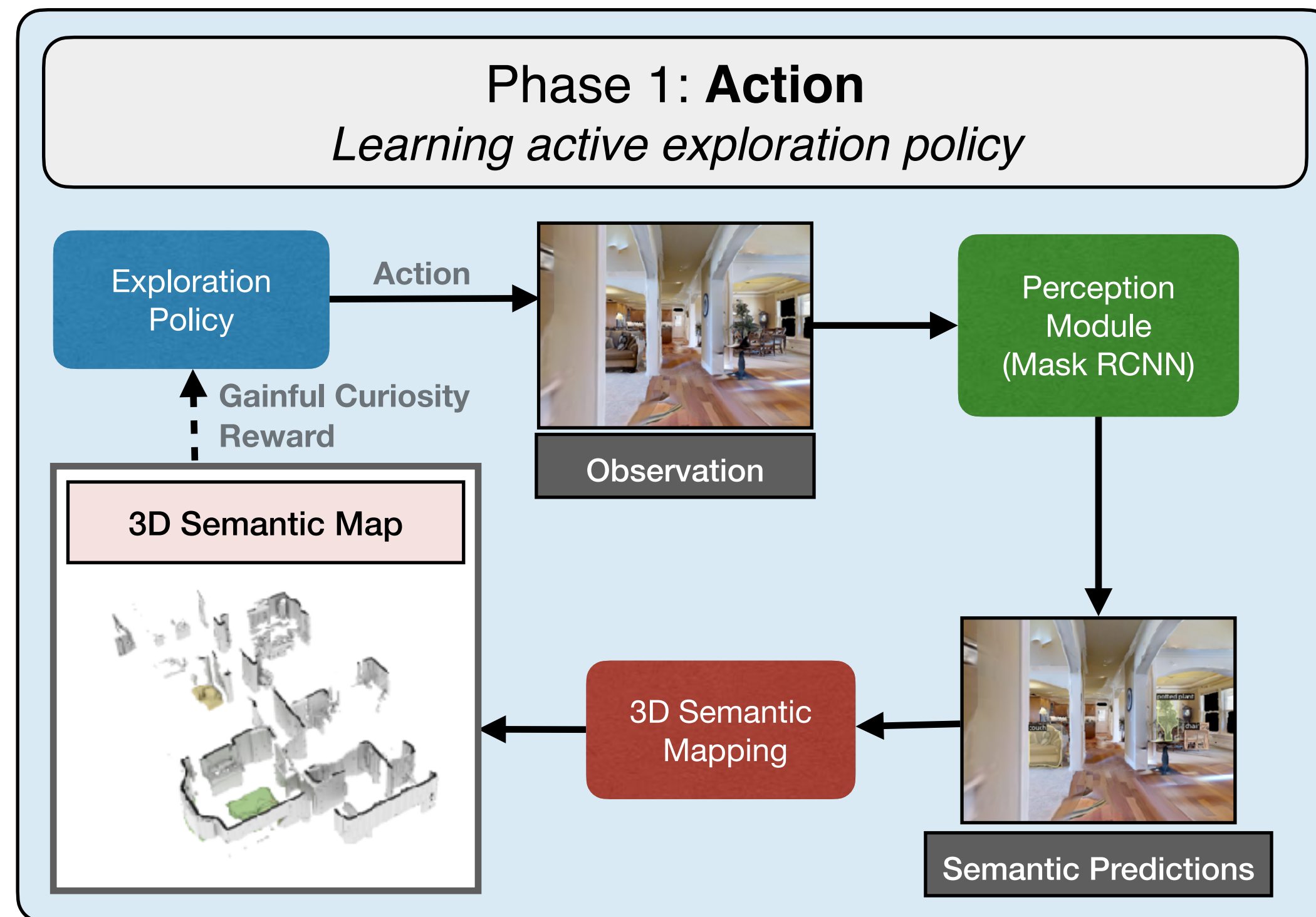


- Chair
- Couch
- Potted Plant
- Bed
- Toilet
- TV

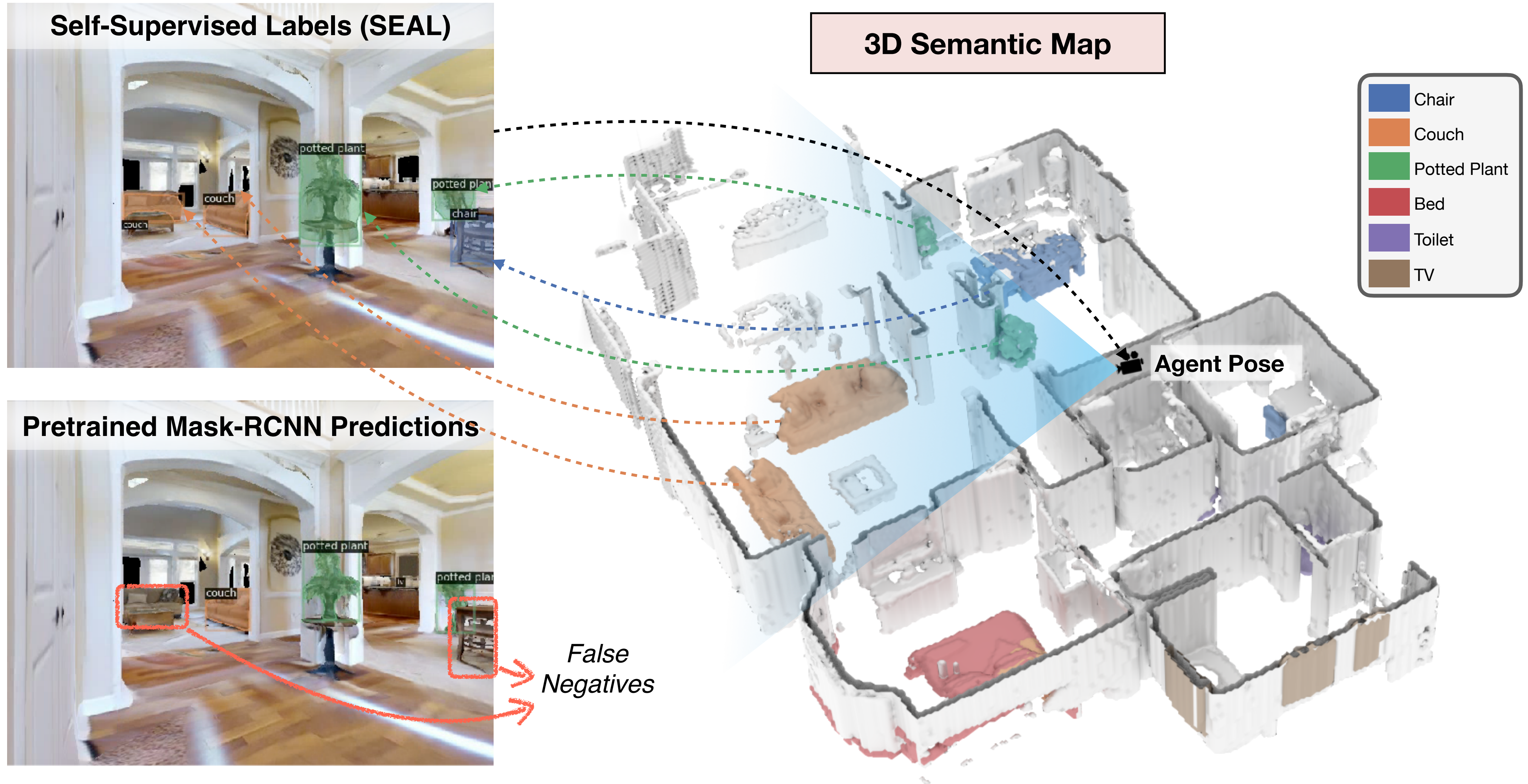
Gainful Curiosity



SEAL: Self-supervised Embodied Active Learning

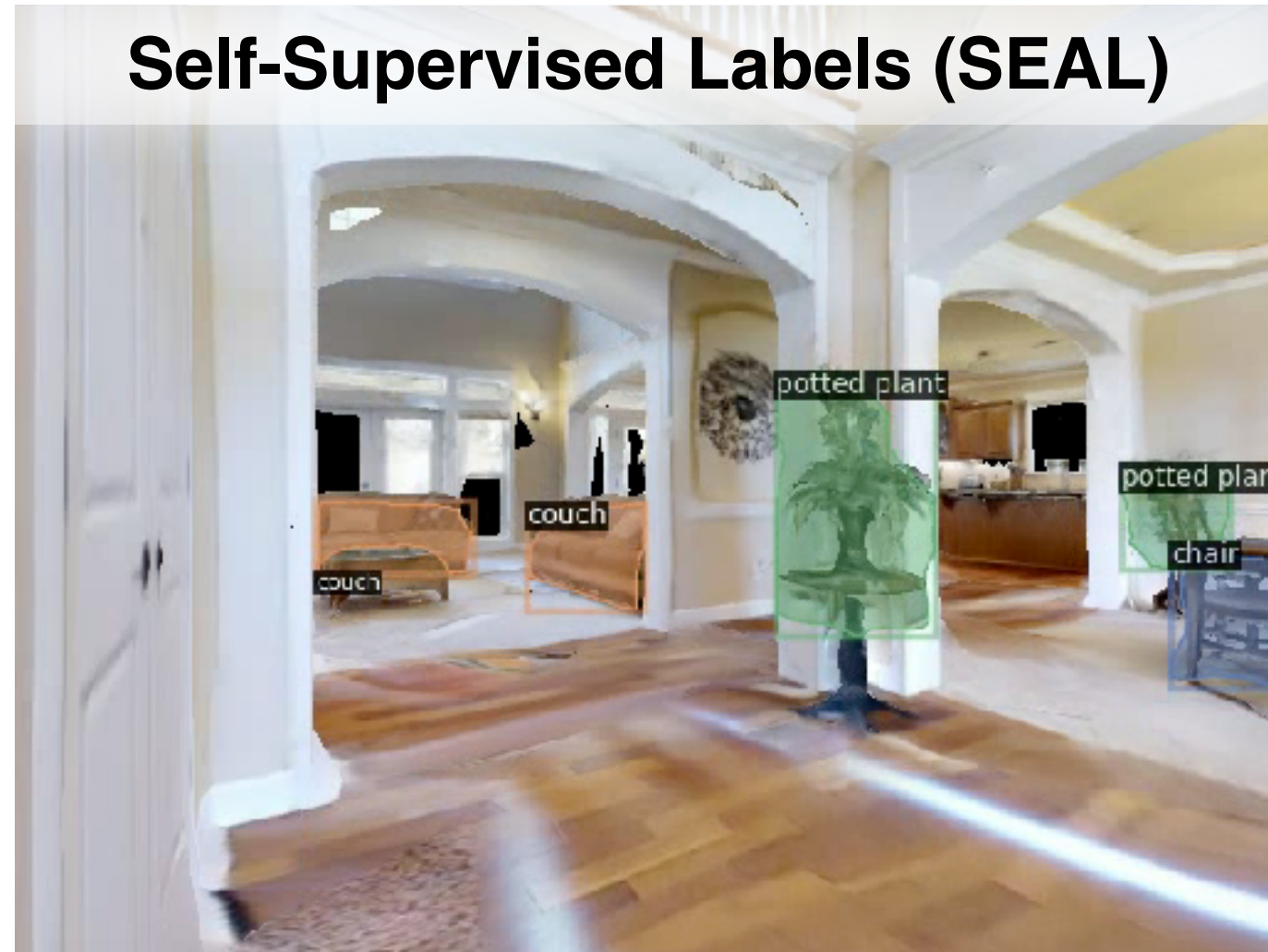


3D Label Propagation



3D Label Propagation

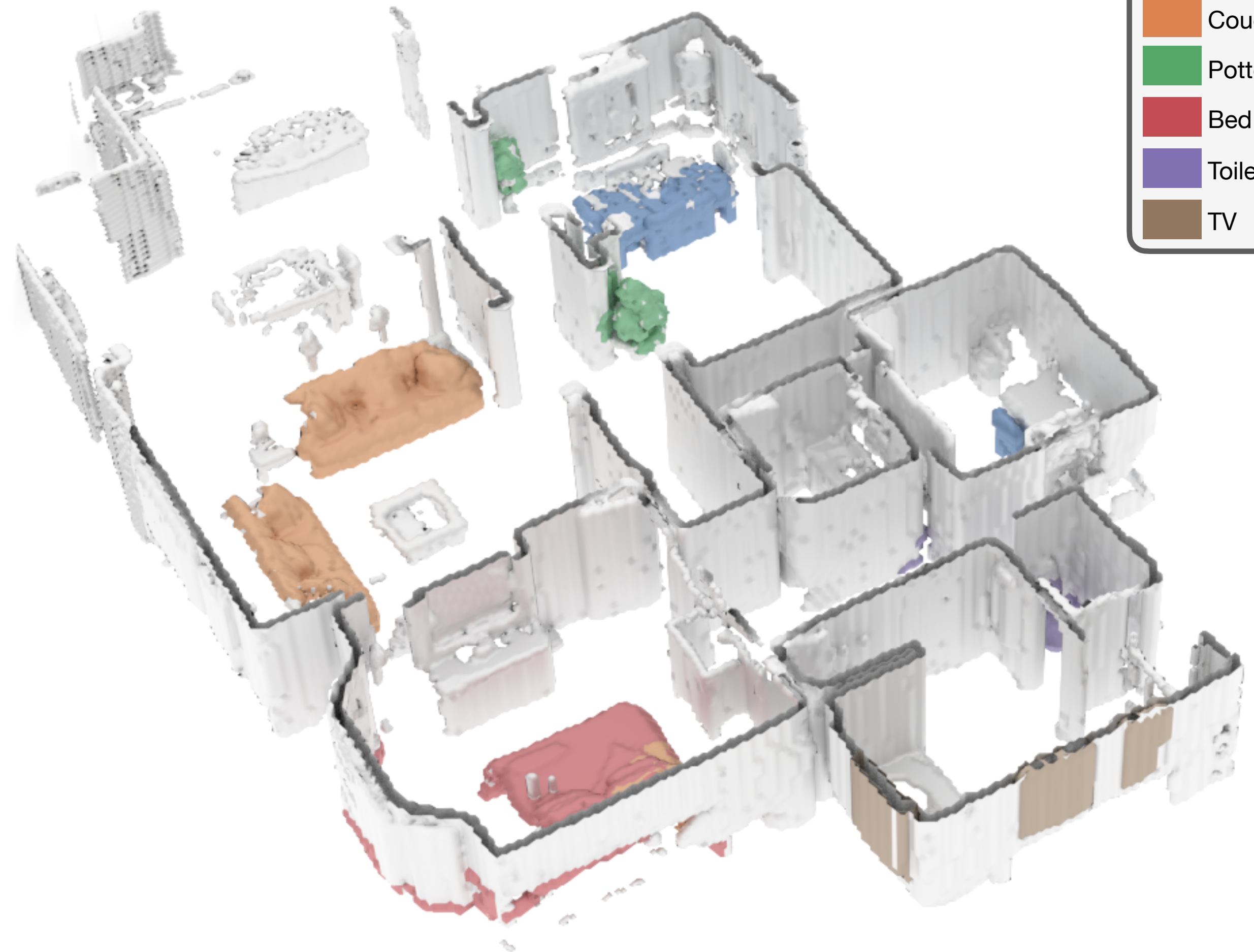
Self-Supervised Labels (SEAL)



Pretrained Mask-RCNN Predictions



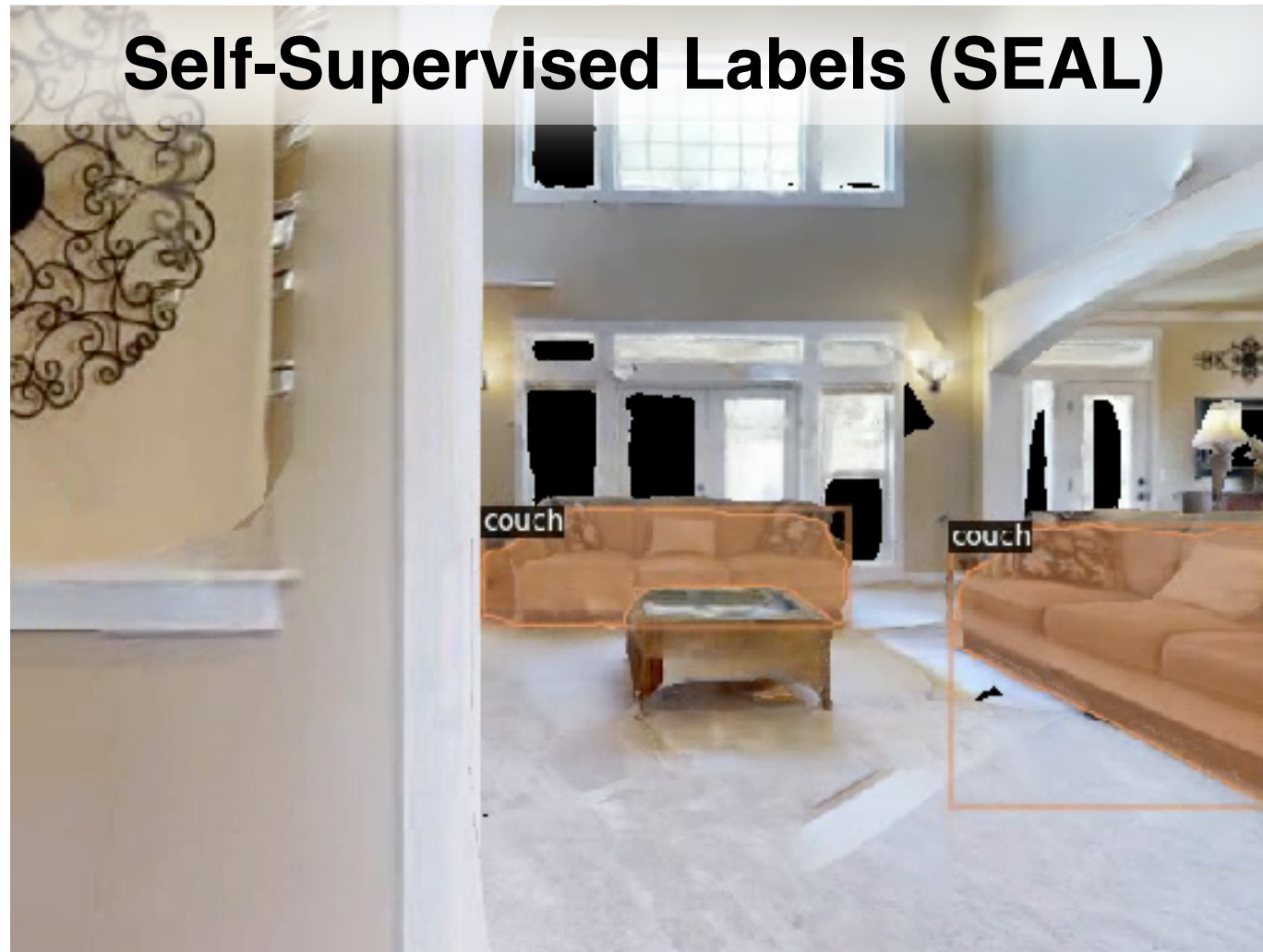
3D Semantic Map



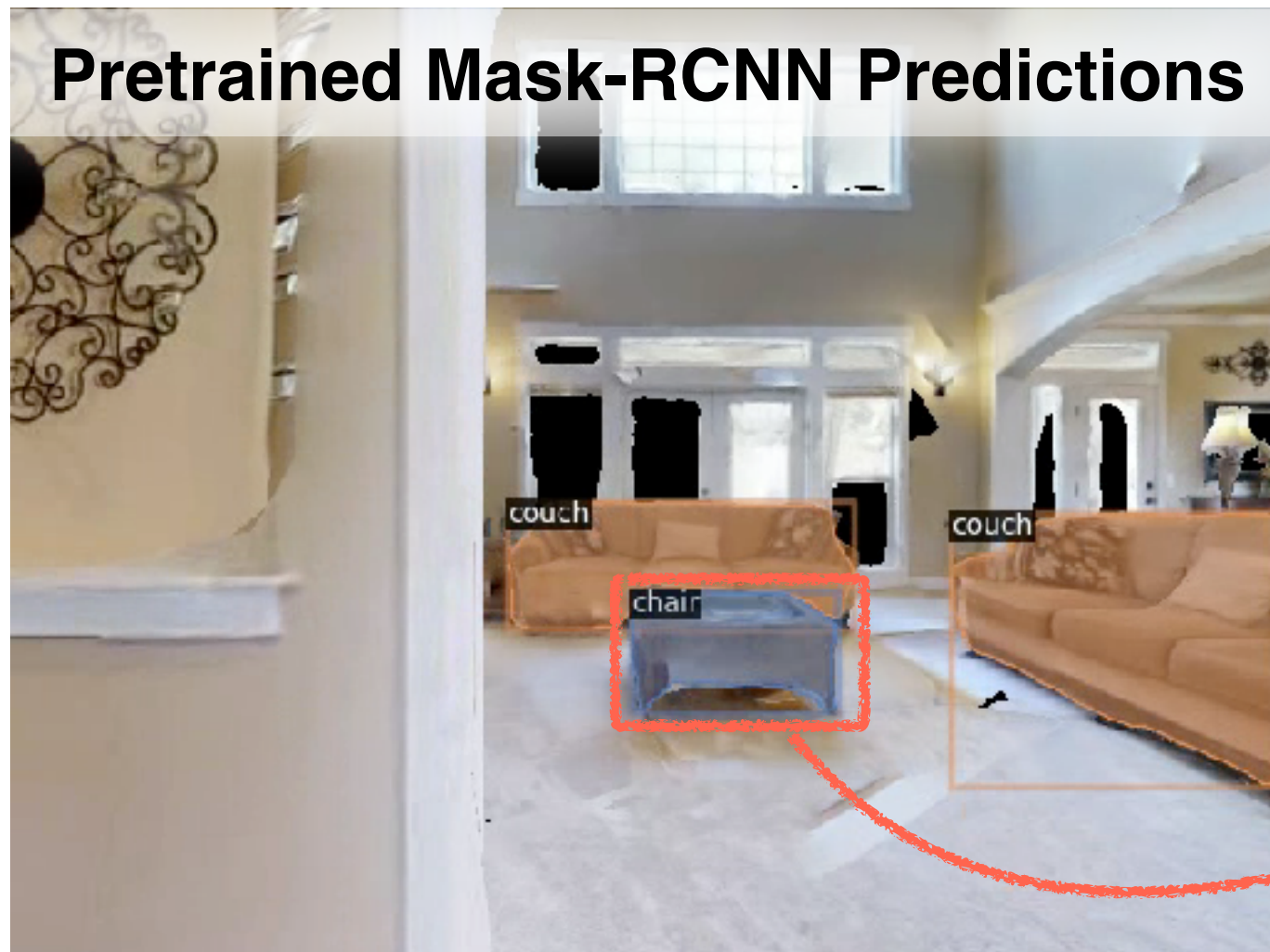
- Chair
- Couch
- Potted Plant
- Bed
- Toilet
- TV

3D Label Propagation

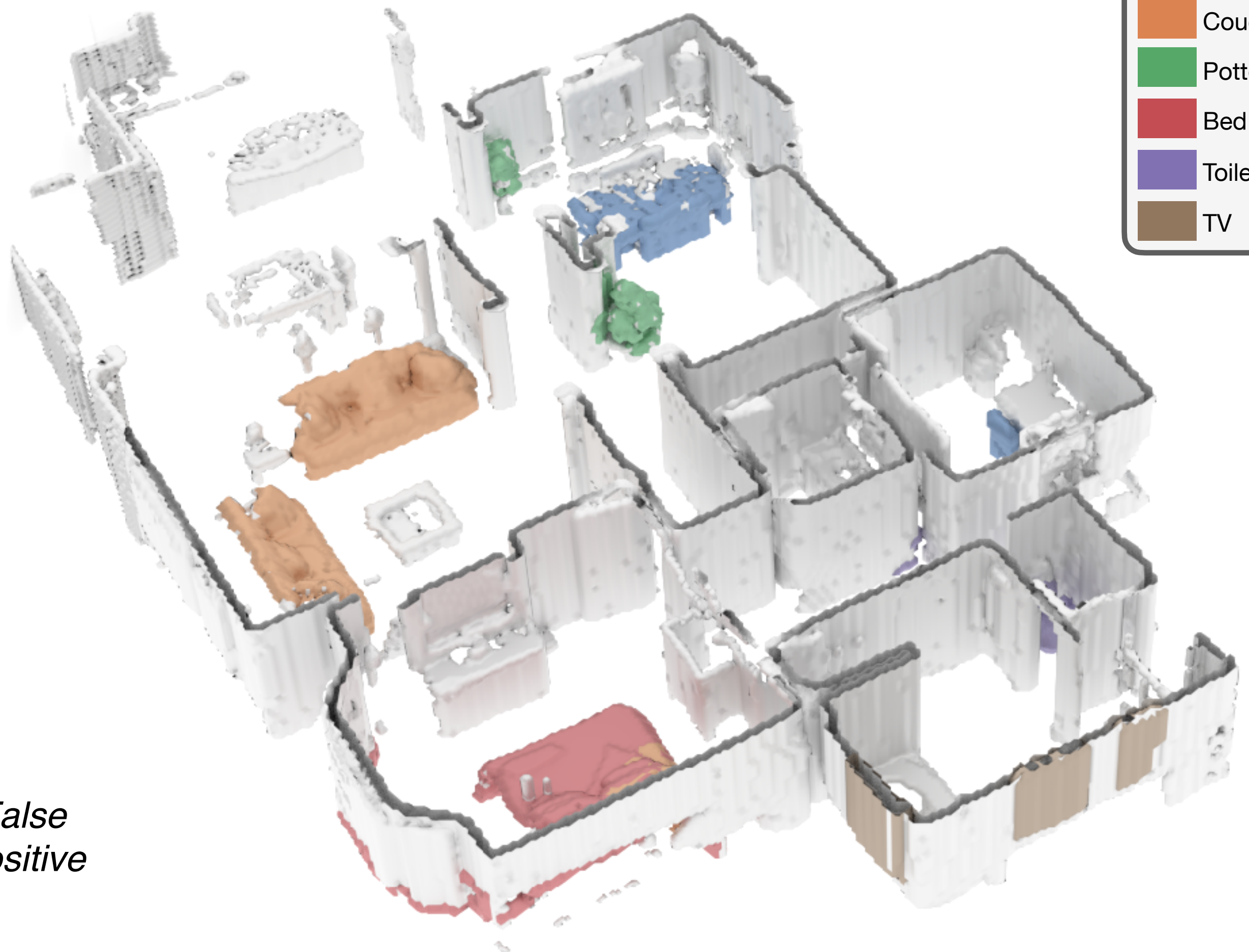
Self-Supervised Labels (SEAL)



Pretrained Mask-RCNN Predictions



3D Semantic Map



- Chair
- Couch
- Potted Plant
- Bed
- Toilet
- TV

False Positive

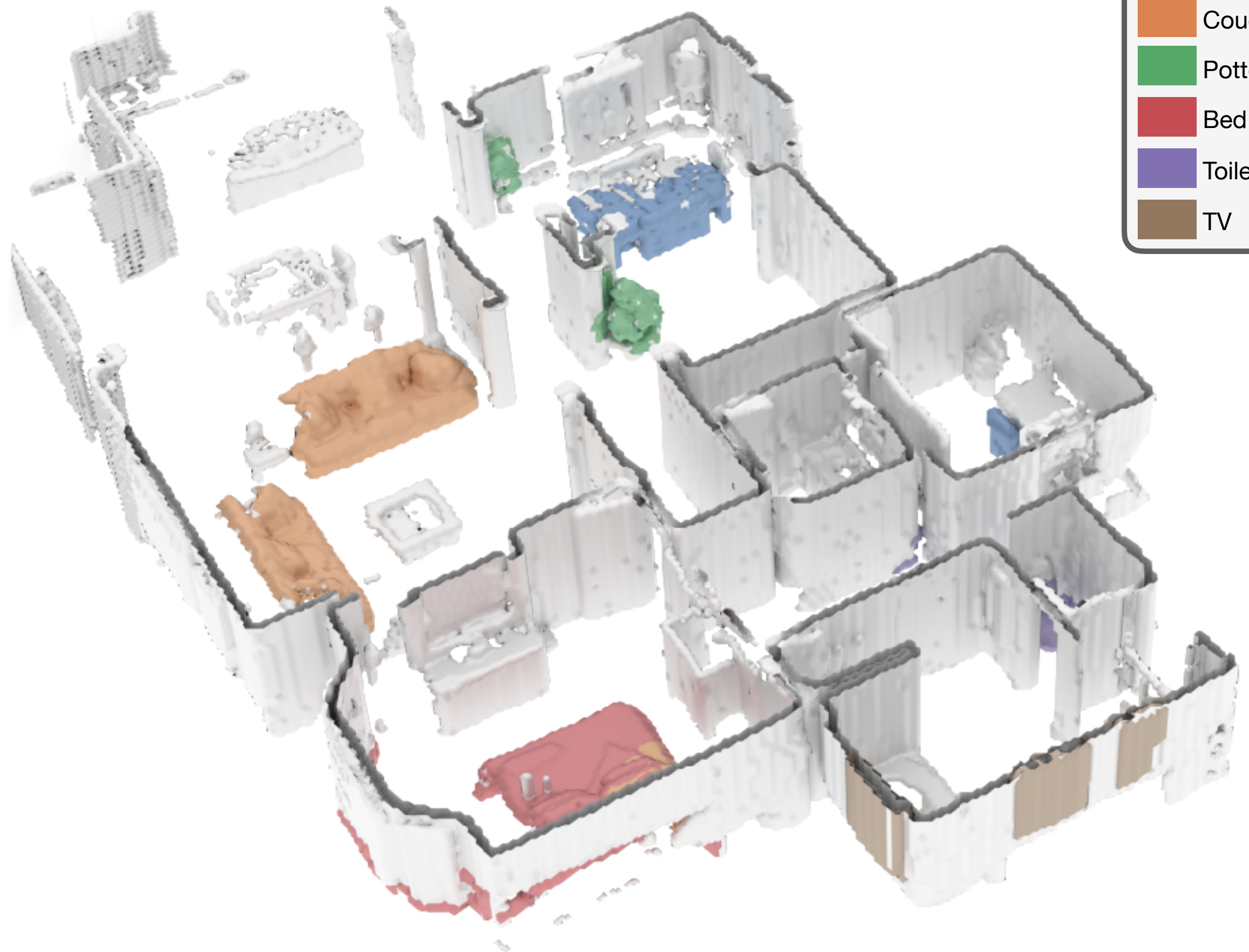
3D Label Propagation



Train
Perception
Model

Perception Model
(Mask RCNN)

3D Semantic Map



Experiments

- Gibson [2] dataset in Habitat simulator [1]
 - 25 train and 5 test scenes
- Training:
 - Action phase: 10 million frames in training scenes
 - Perception phase: 1 trajectory of 300 steps per scene
- Objective: Maximize perception model performance on unseen images in test scenes
 - Metrics:
 - Object Detection **AP50**
 - Instance Segmentation **AP50**

	Action	Perception
Generalization	Train	Train
Specialization	Train	Train + 1 episode test

Results

Metric: AP50

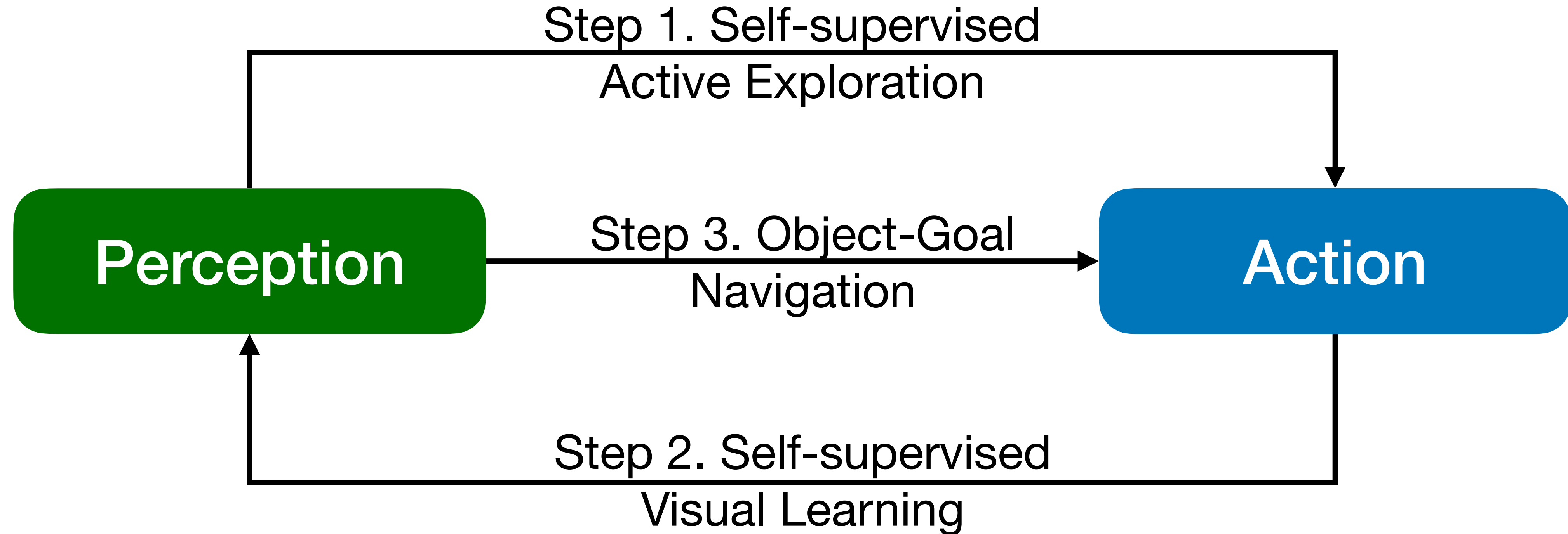
Method	Generalization		Specialization	
	Object Detection	Instance Segmentation	Object Detection	Instance Segmentation
Pretrained Mask-RCNN	34.82	32.54	34.82	32.54
Random Policy + Self-training [1]	33.41	31.89	34.11	31.23
Random Policy + Optical Flow [2]	33.97	32.34	34.33	32.22
Frontier Exploration [3] + Self-training	33.78	32.45	33.29	32.50
Frontier Exploration + Optical Flow	35.22	31.90	34.19	32.12
Active Neural SLAM [4] + Self-training	34.35	31.20	34.84	32.44
Active Neural SLAM + Optical Flow	35.85	32.22	35.90	33.12
Semantic Curiosity [5] + Self-training	35.04	32.19	35.23	32.88
Semantic Curiosity + Optical Flow	35.61	32.57	35.71	33.29
SEAL	40.02	36.23	41.23	37.28

**Self-Supervised
Exploration**

**Single episode
in test scene**

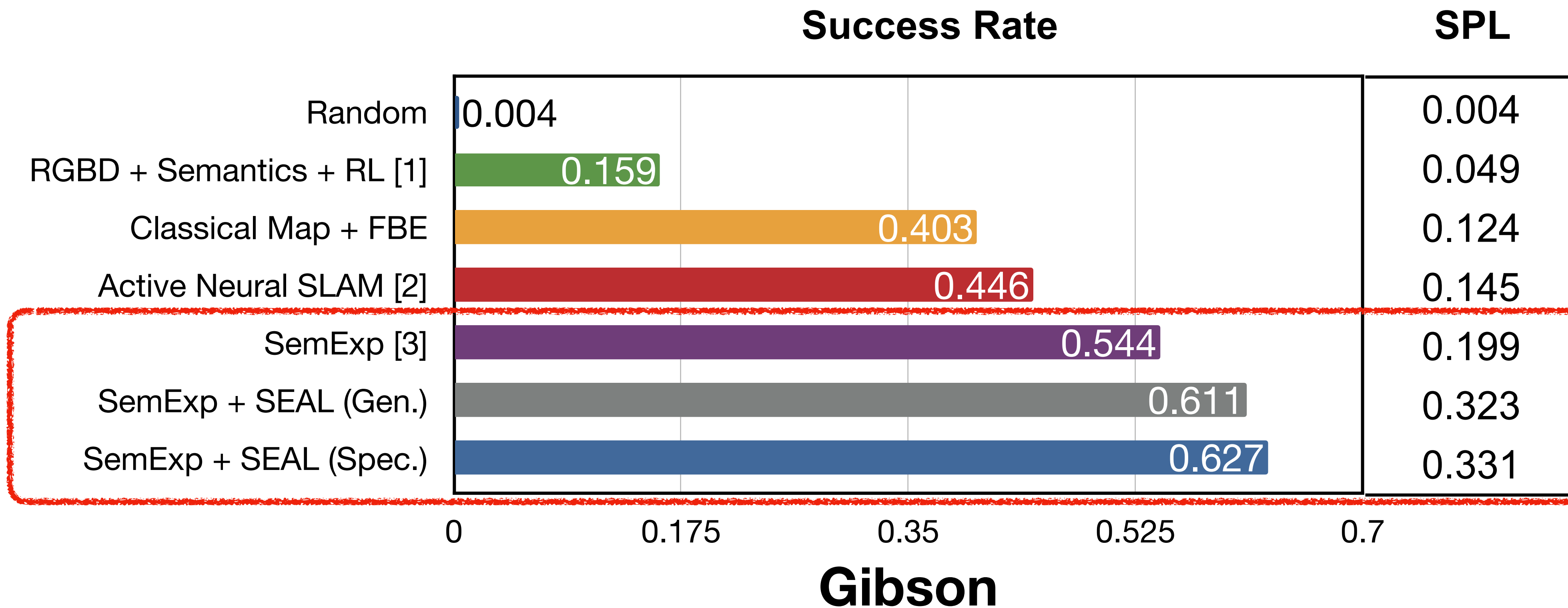
[1] Yalniz et al. 2019, [2] Horn and Schunck. AI 1981, [3] Yamauchi. 1997, [4] Chaplot et al. ICLR 2020, [5] Chaplot et al. ECCV 2020

Perception-Action Loop



We must perceive in order to move, but we must also move in order to perceive
- Gibson (1979)

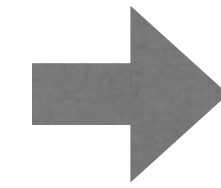
Results: Object Goal Navigation



[1] Mousavian et al. ICRA-19, [2] Chaplot et al. ICLR-20, [3] Chaplot et al. NeurIPS-20

Results: Weak supervision

Mask RCNN



Ground Truth



	Fine-tuning Mask-RCNN		SEAL	
Num labels	Object Detection	Instance Segmentation	Object Detection	Instance Segmentation
0	34.82	32.54	41.23	37.28
5	34.22	31.67	41.44	37.65
10	35.14	32.52	42.63	38.48



SEAL: Self-supervised Embodied Active Learning

Devendra Singh Chaplot, Murtaza Dalal, Saurabh Gupta, Jitendra Malik, Ruslan Salakhutdinov
NeurIPS 2021

Webpage: <https://devendrachaplot.github.io/projects/seal>

Thank you



Devendra Singh Chaplot

Webpage: <http://devendrachaplot.github.io/>

Email: dchaplot@fb.com

Twitter: @dchaplot