

SPIN: Simultaneous Perception, Interaction and Navigation

Anonymous CVPR submission

Paper ID 17

Abstract

001 While there has been remarkable progress recently in
002 the fields of manipulation and locomotion, embodied mo-
003 bile manipulation remains a long-standing challenge. Com-
004 pared to locomotion or static manipulation, a mobile system
005 makes a diverse range of long-horizon tasks feasible in un-
006 structured and dynamic environments. Prior works use dis-
007 entangled modular skills for mobility and manipulation that
008 are trivially tied together, causing several limitations such
009 as compounding errors, delays in decision-making, and no
010 whole-body coordination. We present a reactive mobile
011 manipulation framework that uses an active visual system
012 to consciously perceive and react to its environment using
013 only ego-vision, without any mapping or planning, similar
014 to how humans leverage whole-body and hand-eye coordi-
015 nation. Videos are available at <https://spin-robot.github.io>
016

1. Introduction

017
018 Consider a person trying to carry a coffee cup
019 through clutter. This not only requires naviga-
020 tional planning from start to goal but planning of
021 the whole body to avoid obstacles along the way.
022 Furthermore, due to ego-centric
023 vision, the person needs to ac-
024 tively look around for obstacles.
025 This general form of mobile ma-
026 nipulation necessitates a cou-
027 pled understanding of whole-
028 body control with active percep-
029 tion as a fundamental capability
030 in embodied cognition.



031 The current paradigm tackles this through classical
032 planning-based control which requires apriori knowledge
033 of the precise location of obstacles with a detailed map of
034 the environment. This assumption is impractical in the real
035 world due to computational reasons, and more importantly,
036 because environments are dynamic and keep changing. Hu-
037 mans, on the other hand, do not rely on precise estimates
038 of obstacles and instead use ego-centric vision to navigate

around them in real-time. In an unfamiliar environment,
where to look is informed by where they want to move
(called ‘active perception’), and how they move in return
determines what they can see immediately afterward. This
integrated mobility and perception allows us to see, adapt,
and react to maneuver through cluttered environments.

This paper presents *SPIN*, an end-to-end approach to
Simultaneous Perception, Interaction, and Navigation. We
train a single model using reinforcement learning (RL) that
not only outputs low-level controls for the robot body and
arm but also predicts where should the robot’s ego-centric
camera look at each time step. We evaluate across 6 bench-
marks in simulation and 2 real-world environments outper-
forming the baselines.

2. Method

We want our mobile manipulator to navigate and manip-
ulate objects while avoiding obstacles in clutter. With an
actuated camera with limited FOV (87° horizontal, 58° ver-
tical), it requires one to look around to simultaneously plan
and avoid obstacles. For this challenging problem setup,
we train our robot to navigate inside procedurally generated
clutter in simulation using RL. The robot is only allowed to
perceive part of its environment visible to the camera and,
learns to coordinate its arm, base, and camera motion.

In practice, since training with RL requires many sam-
ples and depth rendering is inefficient, we divide training
into two phases. In the first one, we learn mobile manipula-
tion behaviors via RL using a cheap-to-compute variant of
depth (scandots) and in phase 2 we train a CNN for percep-
tion from depth images as illustrated in Figure 1.

Phase 1 - Learning Simultaneous Perception, Interac- tion and Navigation

In this stage, we use RL to learn to
control all the joints of the robot to navigate clutter and pick
target objects. Since rendering depth images directly from
the robot camera is expensive, we use an ersatz version that
contains the same information and is cheap to compute. We
do so using *scandots* s_t which are the xyz coordinates of the
bounding box of each obstacle. To specify which object to
pick, we give the initial location of the object o_i . Instead of
the object image, we give the current location of the object

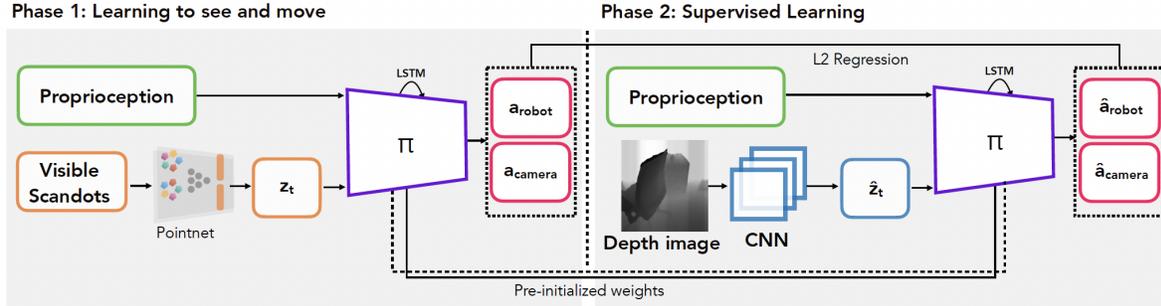


Figure 1. We learn a policy that uses ego-vision to simultaneously perceive, interact, and navigate in clutter. We propose Coupled Visuomotor Optimization (CVO) that learns robot and camera actions at the same time using an RL policy. We only provide scandots if they are visible in the agent’s fov allowing it to learn to move its camera and aggregate information about its environment. This is followed by a phase-2 supervised training where this behavior is distilled into a student network operating with ego-centric depth images.

	Reach						Pick	Place					
	Scenario 1			Scenario 2				Scenario 1			Scenario 2		
	Easy	Medium	Hard	Easy	Medium	Hard		Easy	Medium	Hard	Easy	Medium	Hard
FixCam	1.00	0.53	0.20	1.00	0.50	0.26	0.86	1.00	0.53	0.16	0.97	0.50	0.20
NoPointNet	1.00	0.87	0.57	1.00	0.77	0.63	0.93	1.00	0.83	0.57	1.00	0.77	0.60
Mapping	1.00	1.00	1.00	0.86	1.00	0.97	0.97	1.00	1.00	1.00	1.00	0.90	0.97
SPIN	1.00	0.97	0.93	1.00	1.00	0.93	0.97	1.00	0.97	0.90	1.00	0.97	0.93

Table 1. We report the success rate of each part of the task including reaching (Reach), picking (Pick), and placing (Place) the target object in the desired location. The placing task requires the agent to bring back the object across the obstacles near its start location.

079 \mathbf{o}_t . Here, scandots \mathbf{s}_t and object location \mathbf{o}_t are privileged
080 information which must later be estimated from depth im-
081 ages. Given this, we train two separate LSTM policies π_{nav}
082 and π_{pick} using a dense reward for each of the tasks and early
083 termination for collisions with obstacles.

084 **Phase 2 - From Scandots to Depth** Scandots are not di-
085 rectly observable in the real world and must instead be es-
086 timated from the depth image. We train a convolution net-
087 work C to convert rendered depth images \mathbf{d}_t to percep-
088 tion latents $\tilde{\mathbf{z}}_t$. This latent is passed to a student policy π' to
089 predict the actions $[\tilde{\mathbf{a}}_{\text{robot}}, \tilde{\mathbf{a}}_{\text{cam}}]$. This is supervised using L2
090 loss from the phase 1 actions. The weights for π' are ini-
091 tialized using π . We train this policy using DAgger [3]. For
092 the navigation policy, we optimize

$$093 \min_{C_{\text{nav}}, \pi'_{\text{nav}}} \|\pi'_{\text{nav}}(C_{\text{nav}}(\mathbf{d}_t), \mathbf{x}_t, \mathbf{g}_t) - \pi_{\text{nav}}(\mathbf{z}_t, \mathbf{x}_t, \mathbf{g}_t)\| \quad (1)$$

094 Note that the teacher policy π_{nav} can be trained using either
095 the coupled or decoupled approach. Similarly, for the pick
096 policy, we estimate current object position \mathbf{o}_t from depth

$$097 \min_{C_{\text{pick}}, \pi'_{\text{pick}}} \|\pi'_{\text{pick}}(C_{\text{pick}}(\mathbf{d}_t), \mathbf{x}_t, \mathbf{o}_i) - \pi_{\text{pick}}(\mathbf{z}_t, \mathbf{x}_t, \mathbf{o}_i)\| \quad (2)$$

098 3. Experiments and Results

099 We use Hello Robot [1] for experiments, train our policy us-
100 ing IssacGym [2], and compare against following baselines:

- **FixCam:** Fixed camera without active perception. 101
- **Mapping:** Policy operating on environment map instead 102
of using a moving depth camera. 103
- **NoPointNet:** Using an MLP, instead of a permutation- 104
invariant PointNet architecture for scandots latent. 105

The simulation benchmark has 6 scenes, 2 of each easy, 106
medium, and hard environment. Easy environments have 0- 107
1 obstacles within a 5m goal range. Medium ones have 2-3 108
obstacles within 5m and the hard ones have heavily clut- 109
tered scenes with 5 obstacles within 5m. In each case, Sce- 110
nario 1 comprises a tight 1m wide long corridor which al- 111
lows the agent to not take shortcuts and reach the goal only 112
by navigating through obstacles. Scenario 2 is an L-shaped 113
corridor with a goal at the end. 114

We compare against various baselines as reported in Ta- 115
ble 1. For each scenario, we report the success rate across 116
10 rollouts and 3 seeds. SPIN achieves $\approx 33\%$ higher suc- 117
cess rate than the NoPointNet baseline since permutation in- 118
variant scandots latent makes the optimization problem eas- 119
ier and also generalizes better at test time. SPIN achieves 120
 $\approx 68\%$ higher success rate than the FixCam baseline with 121
the camera pointing straight ahead. SPIN is better than the 122
Mapping baseline because the systematic noise in the object 123
locations makes it hard for the robot to avoid them, espe- 124
cially in cluttered environments, whereas SPIN can contin- 125
uously estimate the position of obstacles while it is moving 126
and adapt the motion online. 127

128 **References**

- 129 [1] Stretch by hello robot. <https://hello-robot.com/>.
130 2
- 131 [2] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo,
132 Michelle Lu, Kier Storey, Miles Macklin, David Hoeller,
133 Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel
134 State. Isaac gym: High performance gpu-based physics simu-
135 lation for robot learning, 2021. 2
- 136 [3] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A re-
137 duction of imitation learning and structured prediction to no-
138 regret online learning. In *Proceedings of the fourteenth in-*
139 *ternational conference on artificial intelligence and statistics*,
140 pages 627–635. JMLR Workshop and Conference Proceed-
141 ings, 2011. 2