

Combining Optimal Control and Learning for Visual Navigation in Novel Environments

Somil Bansal^{*1} Varun Tolani^{*1} Saurabh Gupta² Jitendra Malik^{1,2} Claire Tomlin¹

Abstract: Model-based control is a popular paradigm for robot navigation because it can leverage a known dynamics model to efficiently plan robust robot trajectories. However, it is challenging to use model-based methods in settings where the environment is *a priori* unknown and can only be observed partially through on-board sensors on the robot. In this work, we address this short-coming by coupling model-based control with learning-based perception. The learning-based perception module produces a series of *waypoints* that guide the robot to the goal via a collision-free path. These waypoints are used by a model-based planner to generate a smooth and dynamically feasible trajectory that is executed on the physical system using feedback control. Our experiments in simulated real-world cluttered environments and on an actual ground vehicle demonstrate that the proposed approach can reach goal locations more reliably and efficiently in novel environments as compared to purely geometric mapping-based or end-to-end learning-based alternatives. Our approach does not rely on detailed explicit 3D maps of the environment, works well with low frame rates, and generalizes well from simulation to the real world. Videos describing our approach and experiments are available on the project website³.

1 Introduction

Autonomous robot navigation is a fundamental and well-studied problems in robotics. However, developing a fully autonomous robot that can navigate in *a priori* unknown environments is difficult due to challenges that span dynamics modeling, on-board perception, localization and mapping, trajectory generation, and optimal control.

One way to approach this problem is to generate a globally-consistent geometric map of the environment, and use it to compute a collision-free trajectory to the goal using optimal control and planning schemes. However, the real-time generation of a globally consistent map tends to be computationally expensive, and can be challenging in texture-less environments or in the presence of transparent, shiny objects, or strong ambient lighting [1]. Alternative approaches employ end-to-end learning to side-step this explicit map estimation step. However, such approaches tend to be extremely sample inefficient and highly specialized to the system they were trained on [2].

In this paper, we present a framework for autonomous, vision-based navigation in novel cluttered indoor environments under the assumption of perfect robot state measurement. We take a factorized approach to navigation that uses *learning* to make high-level navigation decisions in unknown environments and leverages *optimal control* to produce smooth trajectories and a robust tracking controller. In particular, we train a Convolutional Neural Network (CNN) that incrementally uses the current RGB image observations to produce a sequence of intermediate states or *waypoints*. These waypoints are produced to guide the robot to the desired target location via a collision-free path in previously unknown environments, and are used as targets for a model-based optimal controller to generate smooth, dynamically feasible control sequences to be executed on the robot. Our approach, *LB-WayPtNav* (Learning-Based WayPoint Navigation), is summarized in Fig. 1.

LB-WayPtNav benefits from the advantages of classical control and learning-based approaches in a way that addresses their individual limitations. Learning can leverage statistical regularities to make

^{*}The first two authors contributed equally to this paper.

¹ University of California, Berkeley.

² Facebook AI Research.

³ Project website: <https://vtolani95.github.io/WayPtNav/>.

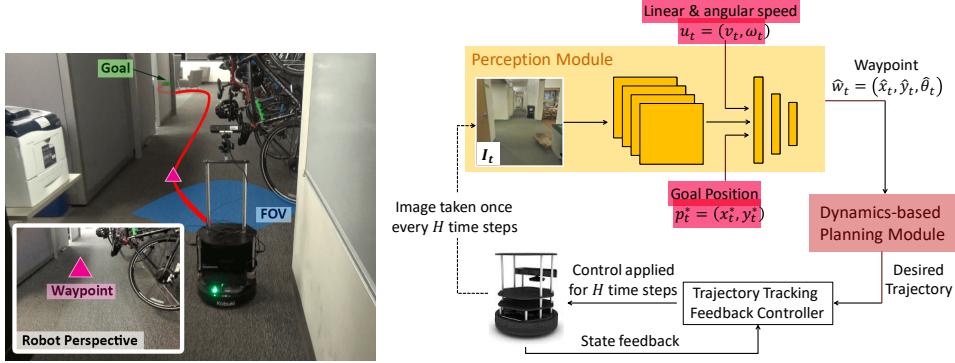


Figure 1. Overview: We consider the problem of navigation from a start position to a goal position. Our approach (LB-WayPtNav) consists of a learning-based perception module and a dynamics model-based planning module. The perception module predicts a waypoint based on the current first-person RGB image observation. This waypoint is used by the model-based planning module to design a controller that smoothly regulates the system to this waypoint. This process is repeated for the next image until the robot reaches the goal.

predictions about the environment from partial views (RGB images) of the environment, allowing generalization to unknown environments. Leveraging underlying dynamics and feedback-based control leads to smooth, continuous, and efficient trajectories that are naturally robust to variations in physical properties and noise in actuation, allowing us to deploy our framework directly from simulation to real-world. Furthermore, learning now does not need to spend interaction samples to learn about the dynamics of the underlying system, and can exclusively focus on dealing with generalization to unknown environments. To summarize, our key contributions are:

- an approach that combines learning and optimal control to robustly maneuver the robot in novel, cluttered environments using only a single on-board RGB camera,
- through simulations and experiments on a mobile robot, we demonstrate that our approach is *better* and more *efficient* at reaching the goals, results in *smoother* trajectories, as compared to End-to-End learning, and more *reliable* than geometric mapping-based approaches,
- we demonstrate that our approach can be *directly* transferred from simulation to unseen, real-world environments without any finetuning or data collection in the real-world,
- an optimal control method for generating optimal waypoints to support large-scale training of deep neural networks for autonomous navigation without requiring any human labeling..

2 Related Work

An extensive body of research studies autonomous navigation. We cannot possibly hope to summarize all these works here, but we attempt to discuss the most closely related approaches.

Classical Robot Navigation. Classical robotics has made significant progress by factorizing the problem of robot navigation into sub-problems of mapping and localization [3, 4], path planning [5], and trajectory tracking. Mapping estimates the 3D structure of the world (using RGB / RGB-D images / LiDAR scans), which is used by a planner to compute paths to goal. However, such purely geometric intermediate representations do not capture navigational affordances (such as: to go far away, one should step into a hallway, *etc.*). Furthermore, mapping is challenging with just RGB observations, and often unreliable even with active depth sensors (such as in presence of shiny or thin objects, or in presence of strong ambient light) [1]. This motivates approaches that leverage object and region semantics during mapping and planning [6, 7]; however, such semantics are often hand-coded. Our work is similarly motivated, but instead of using geometry-based reasoning or hand-crafted heuristics, we employ learning to directly predict good waypoints to convey the robot to desired target locations. This also side-steps the need for explicit map-building.

End-to-End (E2E) Learning for Navigation. There has been a recent interest in employing end-to-end learning for training policies for goal-driven navigation [8, 9, 10, 11]. The typical motivation here is to incorporate semantics and common-sense reasoning into navigation policies. While Zhu *et*

al. [8] learn policies that work well in training environments, Gupta *et al.* [9] and Khan *et al.* [10] design policies that generalize to previously unseen environments. Most such works abstract out dynamics and work with a set of macro-actions (going forward x cm, turning θ°). Such ignorance of dynamics results in jerky and inefficient policies that exhibit stop-and-go behavior on a real robot. Several works also use end-to-end learning for navigation using laser scans [12, 13, 14], for training and combining a local planner with a higher level roadmap for long range navigation [15, 16, 17, 18], or visual servoing [19]. Numerous other works have tackled navigation in synthetic game environments [20, 21, 22], and largely ignore considerations of real-world deployment, such as dynamics and state estimation. Researchers have also employed learning to tackle locomotion problems [23, 24, 25, 26]. These works learn policies for collision-avoidance, *i.e.* how to move around in an environment without colliding. Kahn *et al.* [24] use motion primitives, while Gandhi *et al.* [23], and Sadeghi and Levine [25] use velocity control for locomotion. While all of these works implement policies for collision avoidance via lower level control, our work studies how policy learning itself should be modified for dynamically feasible low-level control for goal-driven behavior.

Combining Optimal Control and Learning. A number of papers seek to combine the best of learning with optimal control for high-speed navigation [27, 28, 29, 30, 31]. Drews *et al.* [28, 32] learn a cost function from monocular images for aggressive race-track driving via Model Predictive Control (MPC). Kaufmann *et al.* [33, 29] use learning to predict waypoints that are used with model-based control for drone racing. The focus of these works is on aggressive control in *training race-track environments*, whereas we seek to learn *goal-driven* policies that work well in *completely novel, cluttered real world testing environments*. This renders their approach for waypoint generation for learning (that does not reason about obstacles explicitly) ineffective. In contrast, waypoints generated for policy learning by our optimal control based method, are guaranteed to generate a collision free trajectory. Muller *et al.* [34] predict waypoints from semantically segmented images and a user provided command for outdoor navigation, and use a PID controller for control. However, they do not explicitly handle obstacles and focuses primarily on lane keeping and making turns. Instead, we use a model-based planner to generate rich, agile, and explicitly dynamically feasible and collision-free control behaviors, without requiring any user commands, and show results in cluttered real-world indoor environments. In a work parallel to ours, Meng *et al.* [35] combine Riemannian Motion Policy with learning for autonomous navigation, whereas we focus on dynamically feasible spline-based policies. Other works such as that from Levine *et al.* [36] and Pan *et al.* [37] combine optimal control and end-to-end learning, by training neural network policies to mimic the optimal control values based on the raw images. We explicitly compare to such an approach in this work.

3 Problem Setup

In this work, we study the problem of autonomous navigation of a ground vehicle in previously unknown indoor environments. We assume that odometry is perfect (*i.e.*, the exact vehicle state is available), and that the environment is static. Dealing with imperfect odometry and dynamic environments are problems in their own right, and we defer them to future work.

We model our ground vehicle as a three-dimensional non-linear Dubins car system with dynamics:

$$\dot{x} = v \cos \phi, \quad \dot{y} = v \sin \phi, \quad \dot{\phi} = \omega, \quad (1)$$

where $z_t := (x_t, y_t, \phi_t)$ is the state of vehicle, $p_t = (x_t, y_t)$ is the position, ϕ_t is the heading, v_t is the speed, and ω_t is the turn rate at time t . The input (control) to the system is $u_t := (v_t, \omega_t)$. The linear and angular speeds v_t and ω_t are bounded within $[0, \bar{v}]$ and $[-\bar{\omega}, \bar{\omega}]$ respectively. We use a discretized version of the dynamics in Eqn. (1) for all planning purposes. The robot is equipped with a monocular RGB camera mounted at a fixed height, oriented at a fixed pitch and pointing forwards. The goal of this paper is to learn control policies for goal-oriented navigation tasks: the robot needs to go to a target position, $p^* = (x^*, y^*)$, specified in the robot's coordinate frame (*e.g.*, 11m forward, 5m left), without colliding with any obstacles. These tasks are to be performed in novel environments whose map or topology is not available to the robot. In particular, at a given time step t , the robot with state z_t receives as input an RGB image of the environment \mathcal{E} , $I_t = I(\mathcal{E}, z_t)$, and the target position $p_t^* = (x_t^*, y_t^*)$ expressed in the current coordinate frame of the robot. The objective is to obtain a control policy that uses these inputs to guide the robot to the target as quickly as possible.

Algorithm 1 Model-based Navigation via Learned Waypoint Prediction

Require: $p^* := (x^*, y^*)$ ▷ Goal location

- 1: **for** $t = 0$ to T **do**
- 2: $z_t := (x_t, y_t, \phi_t)$; $u_t := (v_t, \omega_t)$ ▷ Measured robot pose, and linear and angular speed
- 3: **Every** H **steps do** ▷ Replan every H steps
- 4: $p_t^* := (x_t^*, y_t^*)$ ▷ Goal location in the robot’s coordinate frame
- 5: $\hat{w}_t = \psi(I_t, u_t, p_t^*)$ ▷ Predict next waypoint
- 6: $\{z^*, u^*\}_{t:t+H} = \text{FitSpline}(\hat{w}_t, u_t)$ ▷ Plan spline-based smooth trajectory
- 7: $\{k, K\}_{t:t+H} = LQR(z_{t:t+H}^*, u_{t:t+H}^*)$ ▷ Tracking controller
- 8: $u_{t+1} = K_t(z_t - z_t^*) + k_t$ ▷ Apply control
- 9: **end for**

4 Model-based Learning for Navigation

We use a learning-based waypoint approach to navigation (LB-WayPtNav). The LB-WayPtNav framework is demonstrated in Figure 1 and summarized in Algorithm 1. LB-WayPtNav makes use of two submodules: perception and planning.

4.1 Perception Module

We implement the perception module using a CNN that takes as input a 224×224 pixel RGB image, I_t , captured from the onboard camera, the target position, p_t^* , specified in the vehicle’s current coordinate frame, and vehicle’s current linear and angular speed, u_t , and outputs the desired next state or a waypoint $\hat{w}_t := (\hat{x}_t, \hat{y}_t, \hat{\theta}_t) = \psi(I_t, u_t, p_t^*)$ (Line 5 in Algorithm 1). Intuitively, the network can use the presence of surfaces and furniture objects like floors, tables, and chairs in the scene, alongside the learned priors about their shapes to generate an estimate of the next waypoint, without explicitly building an accurate map of the environment. This allows a more guided and efficient exploration in novel environments based on the robot’s prior experience with similar scenes and objects.

4.2 Planning and Control Module

Given a waypoint \hat{w}_t , and the current linear and angular speed u_t , the planning module uses the system dynamics in Eqn. (1) to design a smooth trajectory, satisfying the dynamics and control constraints, from the current vehicle state to the waypoint. In this work, we represent the x and y trajectories using third-order splines, whose parameters can be obtained using \hat{w}_t and u_t [38]. This corresponds to solving a set of linear equations, and thus, planning can be done efficiently onboard. Since the heading of the vehicle can be determined from the x and y trajectories, a spline-based planner ultimately provides the desired state and control trajectories, $\{z^*, u^*\}_{t:t+H} = \text{FitSpline}(\hat{w}_t, u_t)$, that the robot follows for the time horizon $[t, t + H]$ to reach the waypoint \hat{w}_t (Line 6). Since the splines are third-order, the generated speed and acceleration trajectories are smooth. This is an important consideration for real robots, since jerky trajectories might lead to compounding sensor errors, poor tracking, or hardware damage [39]. While we use splines in this work for computational efficiency, other model-based planning schemes can also be used for trajectory planning.

To track the generated trajectory $\{z^*, u^*\}$, we design a LQR-based linear feedback controller [40], $\{k, K\}_{t:t+H} = LQR(z_{t:t+H}^*, u_{t:t+H}^*)$ (Line 7). Here k and K represent the feed-forward and feedback terms respectively. The LQR controller is obtained using the dynamics in Eqn. (1), linearized around the trajectory $\{z^*, u^*\}$. LQR is a widely used feedback controller in robotic systems to make planning robust to external disturbances and mismatches between the dynamics model and the actual system [41]. This feedback controller allows us to deploy the proposed framework directly from simulation to a real robot (provided the real-world and simulation environments are visually similar), even though the model in Eqn. (1) may not capture the true physics of the robot.

The control commands generated by the LQR controller are executed on the system over a time horizon of H seconds (Line 8), and then a new image is obtained. Consequently, a new waypoint and plan is generated. This entire process is repeated until the robot reaches the goal position.

4.3 Training Details

LB-WayPtNav’s perception module is trained via supervised learning in training environments where the underlying map is known. Even though the environment map is known during training, no such assumption is made during test time and the robot relies only on an RGB image and other on-board sensors. The knowledge of the map during training allows us to compute optimal waypoints (and trajectories) by formulating the navigation problem between randomly sampled pairs of start and goal locations as an optimal control problem, which can be solved using MPC (described in appendix in Section 8.2). The proposed method does not require any human labeling and can be used to generate supervision for a variety of ground and aerial vehicles. Given first-person images and relative goal coordinates as input, the perception module is trained to predict these optimal waypoints.

5 Simulation Experiments

LB-WayPtNav is aimed at combining classical optimal control with learning for interpreting images. In this section, we present experiments in simulation, and compare to representative methods that only use E2E learning (by ignoring all knowledge about the known system), and that only use geometric mapping and path planning (and no learning).

Simulation Setup: Our simulation experiments are conducted in environments derived from scans of real world buildings (from the Stanford large-scale 3D Indoor Spaces dataset [42]). Scans from 2 buildings were used to generate training data to train LB-WayPtNav. 185 test episodes (start, goal position pairs) in a 3rd *held-out* building were used for testing the different methods. Test episodes are sampled to include scenarios such as: going around obstacles, going out of the room, going from one hallway to another. Though training and test environments consists of indoor offices and labs, their layouts and visual appearances are quite different (see Section 8.4 for some images).

Implementation Details: We used a pre-trained ResNet-50 [43] as the CNN backbone for the perception module, and finetuned it for waypoint prediction with MSE loss using the Adam optimizer. More training details and the exact CNN architecture are provided in appendix in Section 8.1.

Comparisons: We compare to two alternative approaches. **E2E Learning:** This approach is trained to directly output the velocity commands corresponding to the optimal trajectories produced by the spline-based planner (the same trajectories used to generate supervision for LB-WayPtNav, see Sec. 8.2). This represents a purely learning based approach that does not explicitly use any system knowledge at test time. **Geometric Mapping and Planning:** This approach represents a learning-free, purely geometric approach. As inferring precise geometry from RGB images is challenging, we provide ideal depth images as input to this approach. These depth images are used to incrementally build up an occupancy map of the environment, that is used with the same spline-based planner (that was used to generate the expert supervision, see Sec. 8.2), to output the velocity controls. Results reported here are with control horizon, $H = 1.5s$. We also tried $H = 0.5, 1.0$, but the results and trends were the same.

Metrics: We make comparisons via the following metrics: success rate (if the robot reaches within $0.3m$ of the goal position without any collisions), the average time to reach the goal (for the successful episodes), and the average acceleration and jerk along the robot trajectory. The latter metrics measure execution smoothness and efficiency with respect to time and power.

5.1 Results

Comparison with the End-to-End learning approach. Table 1 presents quantitative comparisons. We note that LB-WayPtNav conveys the robot to the goal location more often (22% higher success rate), much faster (40% less time to reach the goal), and with less power consumption (50% less acceleration). Figure 2(left) shows top-view visualization of trajectories executed by the two methods. Top-views maps are being only used for visualization, both methods operate purely based on first-person RGB image inputs. As LB-WayPtNav uses a model-based planner to compute exact controls, it only has to learn “where to go” next, as opposed to the E2E method that also needs to learn “how to go” there. Consequently, LB-WayPtNav is able to successfully navigate through narrow hallways, and make tight turns around obstacles and corners, while E2E method struggles. This is further substantiated by the velocity control profiles in Figure 2(right). Even though the E2E method was trained to predict smooth control profiles (as generated by the expert policy), the control profiles at

Table 1. Quantitative Comparisons in Simulation: Various metrics for different approaches across the test navigation tasks: success rate (higher is better), average time to reach goal, jerk and acceleration along the robot trajectory (lower is better) for successful episodes. LB-WayPtNav conveys the robot to the goal location more often, faster, and produces considerably less jerky trajectories than E2E learning approach. Since LB-WayPtNav only uses the current RGB image, whereas the geometric mapping and planning approach integrates information from perfect depth images, it outperforms LB-WayPtNav in simulation. However, performance is comparable when the mapping based approach only uses the current image (like LB-WayPtNav, but still depth vs. RGB).

| Agent | Input | Success (%) | Time taken (s) | Acceleration (m/s^2) | Jerk (m/s^3) |
|----------------------|------------------------|-------------|-------------------|--------------------------|------------------|
| Expert | Full map | 100 | 10.78 ± 2.64 | 0.11 ± 0.03 | 0.36 ± 0.14 |
| LB-WayPtNav (our) | RGB | 80.65 | 11.52 ± 3.00 | 0.10 ± 0.04 | 0.39 ± 0.16 |
| End To End | RGB | 58.06 | 19.16 ± 10.45 | 0.23 ± 0.02 | 8.07 ± 0.94 |
| Mapping (memoryless) | Depth | 86.56 | 10.96 ± 2.74 | 0.11 ± 0.03 | 0.36 ± 0.14 |
| Mapping | Depth + Spatial Memory | 97.85 | 10.95 ± 2.75 | 0.11 ± 0.03 | 0.36 ± 0.14 |

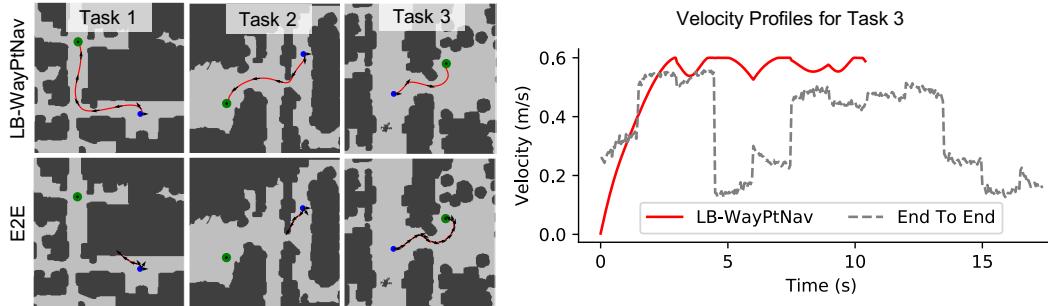


Figure 2. Trajectory Visualization: We visualize the trajectories produced by the model-based planning approach (top row) and the end-to-end (E2E) learning approach (bottom row) for sample test tasks. The E2E learning approach struggles to navigate around the tight corners or narrow hallways, whereas LB-WayPtNav is able to produce a smooth, collision-free trajectory to reach the target position. Even though both approaches are able to reach the target position for task 3, LB-WayPtNav takes only 10s to reach the target whereas the E2E learning approach takes about 17s. Moreover, the control profile produced by the E2E learning approach is significantly more jerky than LB-WayPtNav, which is often concerning for real robots as they are power inefficient, can lead to significant errors in sensors and cause hardware damage.

test time are still discontinuous and jerky. We also experimented with adding a smoothing loss while training the E2E policy; though it helped reduce the average jerk, there was also a significant decline in the success rate. This indicates that learning both an accurate and a smooth control profile can be a hard learning problem. In contrast, as LB-WayPtNav uses model-based control for computing control commands, it achieves average acceleration and jerk that is as low as that of an expert.¹ This distinction has a significant implication for actual robots since the consumed power is directly proportional to the average acceleration. Hence, for the same battery capacity, LB-WayPtNav will drive the robot twice as far as compared to E2E learning.

Comparison with Geometric Mapping and Planning Approach. We note that an online geometric mapping and planning approach, when used with ideal depth image observations, achieves near-expert performance. This is not surprising as perfect depth and egomotion satisfies the exact assumptions made by such an approach. Since LB-WayPtNav is a reactive planning framework, we also compare to a memory-less version that uses a map derived from *only* the current depth image. Even though the performance of memory-less planner is comparable to LB-WayPtNav, it still outperforms slightly due to the perfect depth estimation in simulation. However, since real-world depth sensors are neither perfect nor have an unlimited range, we see a noticeable drop in the performance of mapping-based planners in real-world experiments as discussed in real world experiments in Section 6.

Visualization of Learned Navigation Affordances. We conducted some analysis to understand what cues LB-WayPtNav leverages in order to solve navigation tasks. Figure 3 shows two related navigation tasks where the robot is initialized in the same state, but is tasked to either go inside a close by room (Case A), or to a far away room that is further down the hallway (Case B). LB-WayPtNav

¹For a fair comparison, we report these metrics only over the test tasks that all approaches succeed at.

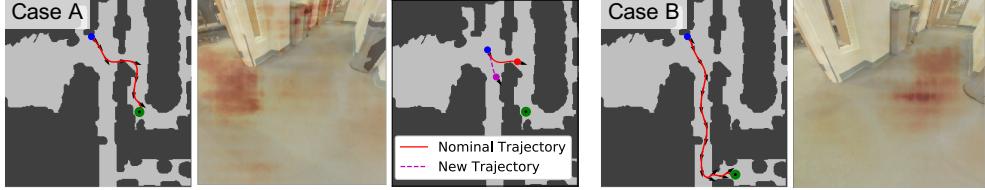


Figure 3. LB-WayPtNav is able to learn the appropriate navigation cues, such as entering the room through the doorway for a goal inside the room, continuing down the hallway for a farther goal. Such cues enable the robot to navigate efficiently in novel environments.

correctly picks appropriate waypoints, and is able to reason that a goal in a far away room is better achieved by going down the hallway as opposed to entering the room currently in front of the robot.

We also conduct an occlusion sensitivity analysis [44], where we measure the change in the predicted waypoint as a rectangular patch is zeroed out at different locations in the input image. We overlay the magnitude of this change in prediction on the input image in Red. LB-WayPtNav focuses on the walls, doorways, hallways and obstacles such as trash-cans as it predicts the next waypoint, and what the network attends to depends on where the robot is trying to go. Furthermore, for Case A, we also show the changed waypoint (in pink) as we zero out the wall pixels. This corresponds to a shorter path to the goal in the absence of the wall. More such examples in the appendix in Section 8.5.

Failure Modes. Even though LB-WayPtNav is able to perform navigation tasks in novel environments, it can only do local reasoning (there is no memory in the network) and gets stuck in some situations. The most prominent failure modes are: a) when the robot is too close to an obstacle, and b) situations that require ‘backtracking’ from an earlier planned path.

6 Hardware Experiments

We next test LB-WayPtNav on a TurtleBot 2 hardware testbed.² We use the network trained in simulation, as described in Section 5, and deploy it directly on the TurtleBot without any additional training or finetuning. We tested the robot in two different buildings, neither of which is in the training dataset (in fact, not even in the S3DIS dataset).³ For state measurement, we use the on-board odometry sensors on the TurtleBot. Test environments for the experiments are described in Table 2.

We repeat each experiment for our method and the three baselines: E2E learning, mapping-based planner, and a memoryless mapping-based planner, for 5 trials each. Results across all 20 trials are summarized in Table 3, where we report success rate, time to reach the goal, acceleration and jerk.

Comparison to E2E learning are consistent with our conclusions from simulation experiments. LB-WayPtNav results in more reliable, faster, and smoother robot trajectories.

Comparison to Geometric Mapping and Planning. Geometric mapping and planning is implemented using the RTAB-Map package [45]. RTAB-Map uses RGB-D images as captured by an on-board RGB-D camera to output an occupancy map that is used with the spline-based planner to output motor commands. As our approach only uses the current image, we also report performance of a memory-less variant of this baseline where occupancy information is derived only from the current observation. While LB-WayPtNav is able to solve 95% of the trials, this memory-less baseline completely fails. It tends to convey the robot too close to obstacles, and fails to recover. In comparison, the map building scheme performs better, with a 40% success rate. This is still a lot lower than performance of our method (95%), and near perfect performance of this scheme in simulation. We investigated the reason for this poor performance, and found that this is largely due to imperfections in depth measurements in the real world. For example, the depth sensor fails to pick-up shiny bike-frames and helmets, black bike-tires and monitor screens, and thin chair legs and power strips on the floor. These systematically missing obstacles cause the robot to collide in experiment 1 and 2. Map quality also substantially deteriorates in the presence of strong ambient light, such as when the sunlight comes in through the window in experiment 4 (visualizations and videos in

² More details about the hardware testbed in Sec. 8.6. Experiment videos are in the supplementary material.

³ Representative images of our experiment environments are shown in Figure 8 in Section 8.4.

Table 2. Experiment setups, with top-views (obtained offline only for visualization), and sample images. Robot starts at the blue dot, and has to arrive at the green dot. Path taken by LB-WayPtNav is shown in red.

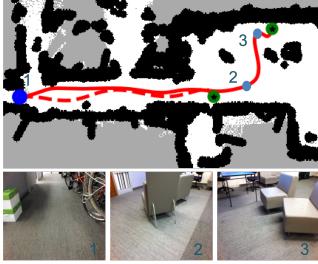
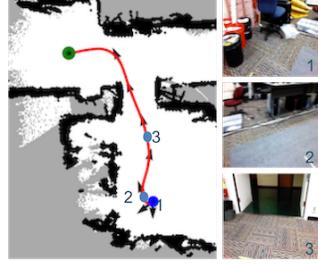
| Experiment 1 and 2 | Experiment 3 | Experiment 4 |
|---|--|---|
|  |  |  |
| Navigation through cluttered environments: This tests if the robot can skillfully pass through clutter in the real world: a narrow hallway with bikes on a bike-rack on one side, and an open space with chairs and sofas. | Leveraging navigation affordances: This tests use of semantic cue for effective navigation. Robot starts inside a room facing a wall. Robot needs to realize it must exit the room through the doorway in order to reach the target location. | Robustness to lighting conditions: Experiment area is similar to that used for experiment 1, but the experiment is performed during the day when sunlight comes from the windows. Robot needs to avoid objects to get to the goal. |

Table 3. Quantitative Comparisons for Hardware Experiments: We deploy LB-WayPtNav and baselines on a TurtleBot 2 hardware testbed for four navigation tasks for 5 trials per task. We report the success rate (higher is better), average time to reach goal, jerk and acceleration along the robot trajectory (lower is better).

| Agent | Input | Success (%) | Time taken (s) | Acceleration (m/s^2) | Jerk (m/s^3) |
|----------------------|------------------------|-------------|------------------|--------------------------|------------------|
| LB-WayPtNav (our) | RGB | 95 | 22.93 ± 2.38 | 0.09 ± 0.01 | 3.01 ± 0.38 |
| End To End | RGB | 50 | 33.88 ± 3.01 | 0.19 ± 0.01 | 6.12 ± 0.18 |
| Mapping (memoryless) | RGB-D | 0 | N/A | N/A | N/A |
| Mapping | RGB-D + Spatial Memory | 40 | 22.13 ± 0.54 | 0.11 ± 0.01 | 3.44 ± 0.21 |

supplementary). These are known fundamental issues with depth sensors, that limit the performance of classical navigation stacks that crucially rely on them.

Performance of LB-WayPtNav: In contrast, our proposed learning-based scheme that leverages robot’s prior experience with similar objects, operates much better without need for extra instrumentation in the form of depth sensors, and without building explicit maps, for the short-horizon tasks that we considered. LB-WayPtNav is able to precisely control the robot through narrow hallways with obstacles (as in experiment 1 and 2) while maintaining a smooth trajectory at all times. This is particularly striking, as the dynamics model used in simulation is only a crude approximation of the physics of a real robot (it does not include any mass and inertia effects, for example). The LQR feedback controller compensates for these approximations, and enables the robot to closely track the desired trajectory (to an accuracy of $4cm$ in experiment 1). LB-WayPtNav also successfully leverages navigation cues (in experiment 3 when it exits the room through a doorway), even when such a behavior was never hard-coded. Furthermore, thanks to the aggressive data augmentation, LB-WayPtNav is able to perform well even under extreme lighting conditions as in Experiment 4.

Furthermore, LB-WayPtNav is agile and reactive. It can adapt to changes in the environment. In an additional experiment (shown in Figure 4), we change the environment as the robot executes its policy. The robot’s goal is to go straight $6m$. It must go around the brown chair. As the policy is executed, we move the chair to repeatedly block the robot’s path (we show the new chair locations in blue and purple, and mark the corresponding positions of the robot at which the chair was moved by same colors). We decrease the control horizon to $0.5s$ for this experiment to allow for a faster visual feedback. The robot successfully reacts to the moving obstacle and reaches the target without colliding.

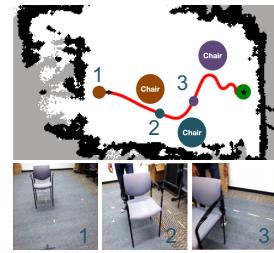


Figure 4. LB-WayPtNav can adapt to dynamic environments.

7 Discussion

We propose LB-WayPtNav, a navigation framework that combines learning and model-based control for goal-driven navigation in novel indoor environments. LB-WayPtNav is better and more reliable at reaching unseen goals compared to an End-to-End learning or a geometric mapping-based approach. Use of a model-based feedback controller allows LB-WayPtNav to successfully generalize from simulation to physical robots.

Even though LB-WayPtNav is somewhat robust to the domain gap between simulation and real-world environments, when the appearances of objects are too different between the two, it fails to predict good waypoints. Thus, it might be desirable to finetune LB-WayPtNav using real-world data. LB-WayPtNav also assumes perfect state estimation and employs a purely reactive policy. These assumptions may not be optimal for long range tasks, wherein incorporating long-term spatial memory in the form of (geometric or learned such as in [9]) maps is often critical. Furthermore, a more detailed study of dynamic environments would be an interesting future work.

Acknowledgments

This research is supported in part by the DARPA Assured Autonomy program under agreement number FA8750-18-C-0101, by NSF under the CPS Frontier project VeHICaL project (1545126), by NSF grants 1739816 and 1837132, by the UC-Philippine-California Advanced Research Institute under project IID-2016-005, by SRC under the CONIX Center, and by Berkeley Deep Drive.

References

- [1] F. Alhwarin, A. Ferrein, and I. Scholl. Ir stereo kinect: improving depth images by combining structured light with ir stereo. In *PRICAI*, 2014.
- [2] B. Recht. A tour of reinforcement learning: The view from continuous control. *Annual Review of Control, Robotics, and Autonomous Systems*, 2018.
- [3] S. Thrun, W. Burgard, and D. Fox. *Probabilistic robotics*. MIT press, 2005.
- [4] J. Fuentes-Pacheco, J. Ruiz-Ascencio, and J. M. Rendón-Mancha. Visual simultaneous localization and mapping: a survey. *Artificial Intelligence Review*, 2015.
- [5] S. M. LaValle. *Planning algorithms*. Cambridge university press, 2006.
- [6] S. L. Bowman, N. Atanasov, K. Daniilidis, and G. J. Pappas. Probabilistic data association for semantic slam. In *ICRA*, 2017.
- [7] B. Kuipers and Y.-T. Byun. A robot exploration and mapping strategy based on a semantic hierarchy of spatial representations. *Robotics and autonomous systems*, 1991.
- [8] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *ICRA*, 2017.
- [9] S. Gupta, V. Tolani, J. Davidson, S. Levine, R. Sukthankar, and J. Malik. Cognitive mapping and planning for visual navigation. *arXiv preprint arXiv:1702.03920*, 2017.
- [10] A. Khan, C. Zhang, N. Atanasov, K. Karydis, V. Kumar, and D. D. Lee. Memory augmented control networks. In *ICLR*, 2018.
- [11] D. K. Kim and T. Chen. Deep neural network for real-time autonomous indoor navigation. *arXiv preprint arXiv:1511.04668*, 2015.
- [12] H.-T. L. Chiang, A. Faust, M. Fiser, and A. Francis. Learning navigation behaviors end-to-end with autorl. *IEEE Robotics and Automation Letters*, 4(2):2007–2014, 2019.
- [13] L. Tai, G. Paolo, and M. Liu. Virtual-to-real deep reinforcement learning: Continuous control of mobile robots for mapless navigation. In *IROS*, 2017.
- [14] W. Zeng, W. Luo, S. Suo, A. Sadat, B. Yang, S. Casas, and R. Urtasun. End-to-end interpretable neural motion planner. In *CVPR*, 2019.
- [15] W. Gao, D. Hsu, W. S. Lee, S. Shen, and K. Subramanian. Intention-net: Integrating planning and deep learning for goal-directed autonomous navigation. In *CoRL*, 2017.
- [16] A. Faust, K. Oslund, O. Ramirez, A. Francis, et al. Prm-rl: Long-range robotic navigation tasks by combining reinforcement learning and sampling-based planning. In *ICRA*, 2018.
- [17] K. Chen, J. P. de Vicente, G. Sepulveda, F. Xia, A. Soto, M. Vazquez, and S. Savarese. A behavioral approach to visual navigation with graph localization networks. In *RSS*, 2019.
- [18] A. Amini, G. Rosman, S. Karaman, and D. Rus. Variational end-to-end navigation and localization. *arXiv preprint arXiv:1811.10119*, 2018.
- [19] F. Sadeghi. Divis: Domain invariant visual servoing for collision-free goal reaching. *RSS*, 2019.
- [20] P. Mirowski, R. Pascanu, F. Viola, H. Soyer, A. Ballard, A. Banino, M. Denil, R. Goroshin, L. Sifre, K. Kavukcuoglu, et al. Learning to navigate in complex environments. In *ICLR*, 2017.
- [21] N. Savinov, A. Dosovitskiy, and V. Koltun. Semi-parametric topological memory for navigation. In *ICLR*, 2018.
- [22] E. Parisotto and R. Salakhutdinov. Neural Map: Structured memory for deep reinforcement learning. In *ICLR*, 2018.
- [23] D. Gandhi, L. Pinto, and A. Gupta. Learning to fly by crashing. In *IROS*, 2017.

- [24] G. Kahn, A. Villaflor, V. Pong, P. Abbeel, and S. Levine. Uncertainty-aware reinforcement learning for collision avoidance. *arXiv preprint arXiv:1702.01182*, 2017.
- [25] F. Sadeghi and S. Levine. (CAD)²RL: Real single-image flight without a single real image. In *RSS*, 2017.
- [26] K. Kang, S. Belkhale, G. Kahn, P. Abbeel, and S. Levine. Generalization through simulation: Integrating simulated and real data into deep reinforcement learning for vision-based autonomous flight. *arXiv preprint arXiv:1902.03701*, 2019.
- [27] C. Richter, W. Vega-Brown, and N. Roy. Bayesian learning for safe high-speed navigation in unknown environments. In *Robotics Research*, pages 325–341. Springer, 2018.
- [28] P. Drews, G. Williams, B. Goldfain, E. A. Theodorou, and J. M. Rehg. Aggressive deep driving: Combining convolutional neural networks and model predictive control. In *CoRL*, 2017.
- [29] E. Kaufmann, A. Loquercio, R. Ranftl, A. Dosovitskiy, V. Koltun, and D. Scaramuzza. Deep drone racing: Learning agile flight in dynamic environments. In *CoRL*, 2018.
- [30] S. Jung, S. Hwang, H. Shin, and D. H. Shim. Perception, guidance, and navigation for indoor autonomous drone racing using deep learning. *IEEE Robotics and Automation Letters*, 2018.
- [31] A. Loquercio, A. I. Maqueda, C. R. del Blanco, and D. Scaramuzza. Dronet: Learning to fly by driving. *IEEE Robotics and Automation Letters*, 3(2):1088–1095, 2018.
- [32] P. Drews, G. Williams, B. Goldfain, E. A. Theodorou, and J. M. Rehg. Vision-based high-speed driving with a deep dynamic observer. *IEEE Robotics and Automation Letters*, 2019.
- [33] E. Kaufmann, M. Gehrig, P. Foehn, R. Ranftl, A. Dosovitskiy, V. Koltun, and D. Scaramuzza. Beauty and the beast: Optimal methods meet learning for drone racing. In *ICRA*, 2019.
- [34] M. Müller, A. Dosovitskiy, B. Ghanem, and V. Koltun. Driving policy transfer via modularity and abstraction. *arXiv preprint arXiv:1804.09364*, 2018.
- [35] X. Meng, N. Ratliff, Y. Xiang, and D. Fox. Neural autonomous navigation with riemannian motion policy. *arXiv preprint arXiv:1904.01762*, 2019.
- [36] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies. *JMLR*, 2016.
- [37] Y. Pan, C.-A. Cheng, K. Saigol, K. Lee, X. Yan, E. Theodorou, and B. Boots. Agile off-road autonomous driving using end-to-end deep imitation learning. In *RSS*, 2018.
- [38] R. Walambe, N. Agarwal, S. Kale, and V. Joshi. Optimal trajectory generation for car-type mobile robot using spline interpolation. *IFAC*, 49(1):601 – 606, 2016.
- [39] D. González, J. Pérez, V. Milanés, and F. Nashashibi. A review of motion planning techniques for automated vehicles. *IEEE Trans. Intelligent Transportation Systems*, 2016.
- [40] D. J. Bender and A. J. Laub. The linear-quadratic optimal regulator for descriptor systems: discrete-time case. *Automatica*, 23(1):71–85, 1987.
- [41] J. Van Den Berg, P. Abbeel, and K. Goldberg. LQG-MP: Optimized path planning for robots with motion uncertainty and imperfect state information. *IJRR*, 2011.
- [42] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese. 3d semantic parsing of large-scale indoor spaces. In *CVPR*, 2016.
- [43] K. He, X. Zhang, S. Ren, J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [44] M. Zeiler, R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014.
- [45] M. Labb   and F. Michaud. Rtab-map as an open-source lidar and visual simultaneous localization and mapping library for large-scale and long-term online operation. *JFR*, 2019.

- [46] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, 2016.
- [47] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *JMLR*, 2014.
- [48] P. Y. Simard, D. Steinkraus, J. C. Platt, et al. Best practices for convolutional neural networks applied to visual document analysis. In *ICDAR*, 2003.
- [49] T. J. Koo and S. Sastry. Differential flatness based full authority helicopter control design. In *CDC*, 1999.
- [50] D. Mellinger and V. Kumar. Minimum snap trajectory generation and control for quadrotors. In *ICRA*, 2011.

8 Supplementary Material

8.1 Network Architecture and Training Details

Implementation Details: We train LB-WayPtNav and E2E agents on 125K data points generated by our expert policy (Section 8.2). All our models are trained with a single GPU worker using TensorFlow [46]. We use MSE loss on the waypoint prediction (respectively on the control command prediction) for training the CNN in our perception module (respectively for E2E learning). We use Adam optimizer to optimize the loss function with a batch size of 64. We train both networks for 35K iterations with a constant learning rate of 10^{-4} and use a weight decay of 10^{-6} to regularize the network. We use ResNet-50 [43], pre-trained for ImageNet Classification, to encode the image input. We remove the top layer, and use a downsampling convolution layer, followed by 5 fully connected layers with 128 neurons each to regress to the optimal waypoint (or control commands for E2E learning). The image features obtained at the last convolution layer are concatenated with the egocentric target position and the current linear and angular speed before passing them to the fully connected layers (see Figure 1). During training, the ResNet layers are finetuned along with the fully connected layers to learn the features that are relevant to the navigation tasks. We use standard techniques used to avoid overfitting including dropout following each fully connected layer except the last (with dropout probability 20%) [47], and data augmentation such as randomly distorting brightness, contrast, adding blur, perspective distortion at training time [48]. Adding these distortions significantly improves the generalization capability of our framework to unseen environments. More details about the kinds of distortions and their effect on the performance can be found in Section 8.3.

8.2 Expert supervision

To generate supervision for training the perception network, we use a Model Predictive Control (MPC) scheme to find a sequence of dynamically feasible waypoints and the corresponding spline trajectories that navigate the robot from the starting state to the goal state. This can be done during the training phase because a map of the environment is available during the training time. However, no such privileged information is used during the test time.

To generate the expert trajectory, we define a cost function that trades-off the distance from the target position and the nearest obstacle, and optimize for the waypoint such that the resultant spline trajectory to the waypoint minimizes this cost function. More specifically, given the map of the environment, we compute the signed distance function to the obstacles, $d^{obs}(x, y)$, at any given position (x, y) . Given the goal position, p^* , of the vehicle, we compute the minimum collision-free distance to the goal, $d^{goal}(x, y)$ (also known as the FMM distance to the goal). The overall cost function to the MPC problem is then given by:

$$J(\mathbf{z}, \mathbf{u}) = \sum_{i=0}^T J_i(z_i, u_i), \quad (2)$$

$$J_i(z_i, u_i) := (\max\{0, \lambda_1 - d^{obs}(x_i, y_i)\})^3 + \lambda_2 (d^{goal}(x_i, y_i))^2, \quad (3)$$

where $\mathbf{z} := (z_0, z_1, \dots, z_T)$ is the state trajectory, \mathbf{u} is the corresponding control trajectory, T is the maximum time horizon, and λ_2 trades-off the distance from the goal position and the obstacles. We only penalize for the obstacle cost when the corresponding robot trajectory is within a distance of λ_1 to an obstacle. Moreover, the obstacle cost is penalized more heavily compared to the goal distance (a cubic penalty vs a quadratic penalty). This is done to ensure that the vehicle trajectory does not go too close to the obstacles. We empirically found that it is significantly harder to learn to accurately predict the waypoints when the vehicle trajectory goes too close to the obstacles, as the prediction error margin for the neural network is much smaller in such a case, leading to a much higher collision rate. Given the cost function in (2), the overall MPC problem is given as

$$\min_{\mathbf{z}, \mathbf{u}} J(\mathbf{z}, \mathbf{u}) \quad (4)$$

$$\text{subject to } x_{i+1} = x_i + \Delta T v_i \cos \phi_i, \quad y_{i+1} = y_i + v_i \sin \phi_i, \quad \phi_{i+1} = \phi_i + \Delta T \omega_i, \quad (5)$$

$$v_i \in [0, \bar{v}], \quad \omega_i \in [-\bar{\omega}, \bar{\omega}], \quad (6)$$

$$z_0 = (0, 0, 0), \quad u_0 = (0, 0), \quad (7)$$

where ΔT is the discretization step for the dynamics in (1), and the initial state and speed are assumed to be zero without loss of generality.

Starting from $i = 0$, we solve the MPC problem in (4) in a receding horizon fashion. In particular, at any timestep $i = t$, we find a waypoint such that the corresponding spline trajectory respects the dynamics and the control constraints in (5) and (6) (and hence is a dynamically feasible trajectory), and minimizes the cost in (4) over a time horizon of H_1 . Thus, we solve the following optimization problem:

$$\min_{\hat{w}_t} \sum_{i=t}^{t+H_1} J_i(z_i, u_i) \quad (8)$$

$$\text{subject to } \{z, u\}_{t:t+H_1} = \text{FitSpline}(\hat{w}_t, u_t), \quad (9)$$

$$z_t, u_t - \text{Given} \quad (10)$$

where $\hat{w}_t := (\hat{x}_t, \hat{y}_t, \hat{\theta}_t)$ is the waypoint, and $\{z, u\}_{t:t+H_1}$ are the corresponding state and control spline trajectories that satisfy the dynamics and the control constraints in (5) and (6), and respect the boundary conditions on the trajectories imposed by z_t , u_t and \hat{w}_t . Such spline trajectories can be computed for a variety of aerial and ground vehicles (see [49, 50] for more details on the aerial vehicles and [38] for the ground vehicles). Thus, the feasible waypoints must be reachable from the current state and speed of the vehicle while respecting the vehicle's control bounds. For example, a pure rotation waypoint (i.e., $\hat{w}_t = (0, 0, \hat{\theta})$) is not feasible for the system if it has a non-zero linear speed at time t , and thus will not be considered as a candidate solution of (8). These dynamics considerations are often ignored in the learning-based methods in literature which work with a set of macro-actions. This often results in an undesirable, jerky, and stop-and-go behavior on a real robot. In this work, we use third-order polynomial splines, whose coefficients are computed using the values of z_t , u_t and \hat{w}_t (see [38] for more details). Use of third-order splines make sure that the velocity and acceleration profiles of the vehicle are smooth, and respect the control bounds.

Let the optimal waypoint corresponding to the optimization problem in (8) be \hat{w}_t^* , and the corresponding optimal trajectories be $\{z^*, u^*\}_{t:t+H_1}$. The image obtained at state z_t^* , I_t , the relative goal position p_t^* , the speed of the robot u_t^* , and the optimal waypoint \hat{w}_t^* thus constitutes one data point for the training. For the End-to-End learning, we use $u_{t:t+H}^*$ instead of \hat{w}_t^* for training the network.

Given the low dimension of the waypoint (three in our case), we use a sampling-based approach to compute the optimal waypoint in (8). We sample the waypoints within the ground-projected field-of-view of the vehicle (ground projected assuming no obstacles). Note, even though we use a sampling-based approach to obtain the optimal waypoint, other gradient-based optimization schemes can also be used, especially when the state space of the vehicle is high-dimensional. We next apply the control sequence u^* for the time horizon $[t, t + H]$ to obtain the state z_{t+H}^* , and repeat the entire procedure in (8) starting from time $t + H$. We continue this process until the robot reaches the goal position. In our work, we use $\lambda_1 = 0.3m$, $\Delta T = 0.05s$, $H_1 = 6s$ and $H = 1.5s$.

The procedure outlined in this section allows us to compute large training datasets using optimal control without requiring any explicit human labeling. Moreover, the generated waypoints are guaranteed to satisfy the dynamics and control constraints. Finally, the cost function in (2) allows us to explicitly ensure that the robot trajectory to the waypoint is collision-free, which is crucial in cluttered indoor environments. It is also worthwhile to note that this procedure can be applied to a variety of ground vehicles and aerial vehicles that are differentially flat. For differentially flat systems, the path planning can be done with respect to a much lower-dimensional state space (or waypoint) using spline trajectories [49, 38, 50], which makes the path planning tractable in real-time.

It is also important to note that intuitively the waypoint in our framework summarizes the local obstacle information in the environment, and hence, its representation choice is very crucial. For example, the choice of $\hat{\theta}$ in \hat{w} affects the shape of the trajectory the vehicle takes to the waypoint (see Figure 5). Typically, $\hat{\theta}$ is chosen as the line of sight angle between the robot's current position and the desired position (\hat{x}, \hat{y}) [29, 34]. However, $\hat{\theta}$ is an extra degree-of-freedom that can be chosen appropriately to generate rich, agile and collision-free trajectories in cluttered environments.

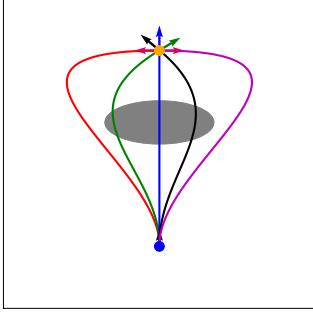


Figure 5. We visualize the spline trajectories produced by our planner for different waypoint angles ($\hat{\theta}$), starting from the same initial state, initial speed and to the same waypoint position. The gray region denotes an obstacle. Due to the dynamics constraints, $\hat{\theta}$ significantly affects the shape of the vehicle trajectory, and hence needs to be chosen appropriately to obtain a collision-free trajectory. In particular, the line of sight trajectory to the waypoint (the Blue trajectory) leads to a collision in this case, whereas if $\hat{\theta}$ is chosen appropriately, a smooth, agile trajectory that goes around the obstacle can be obtained.

8.3 Image distortions and domain randomization

During training, we apply a variety of random distortions to images, such as adding blur, removing some pixels, adding some superpixels, changing image saturation, sharpness, contrast and brightness. We also apply perspective distortions to images, such as varying the field-of-view and tilt (pitch) of the camera. For adding these distortions, we use the *imgaug* python library.

Some examples of sample distortions are shown in Figure 6. Adding these distortions during training significantly improves the generalization of the trained network to unseen environments. For example, for LB-WayPtNav, adding distortions increases the success rate from 47.94% to 80.65%. Adding perspective distortions particularly improves the generalization to the real-world systems, for which the camera tilt will inevitably change as the robot is moving through the environment.

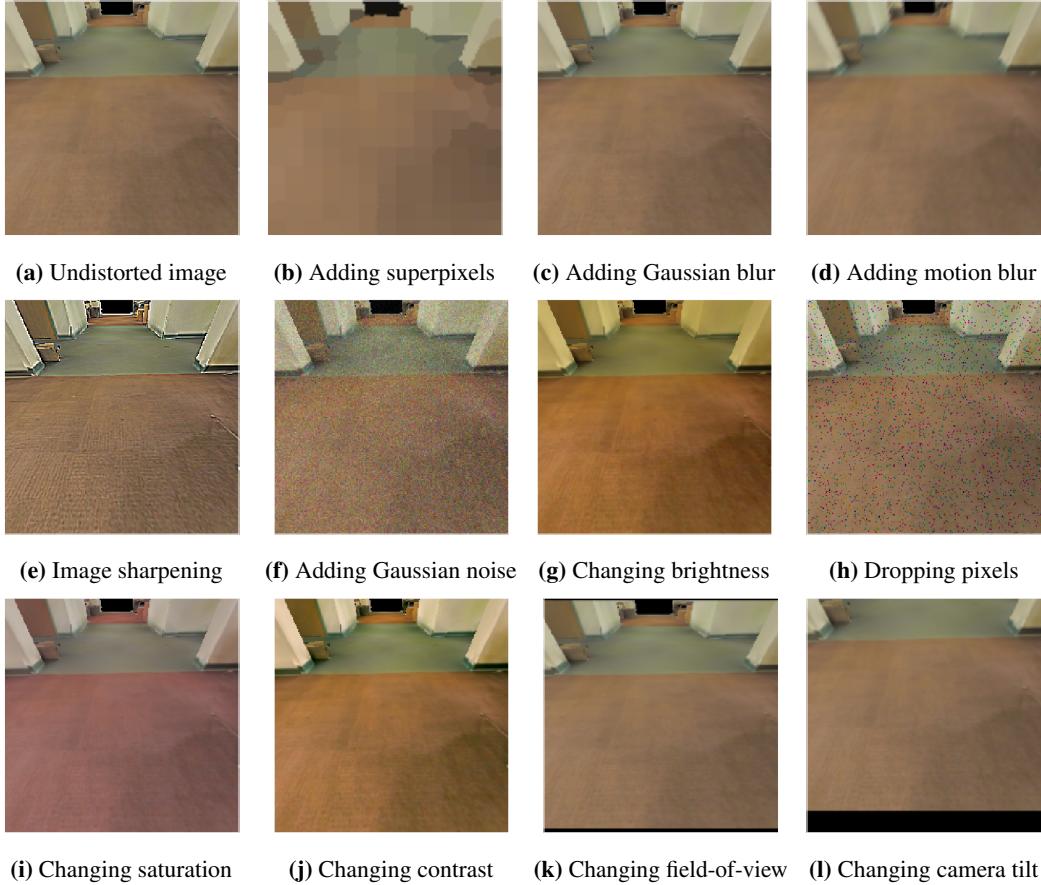


Figure 6. Examples of several image distortions that have been randomly applied during the training phase. The actual undistorted image is shown in (a). Adding these distortions significantly improves the generalization capability of our framework to unseen environments.

8.4 Training and test areas

We illustrate some representative scenes from the test and the training buildings from the S3DIS dataset in Figure 7. Note that the navigation tasks were not limited to these scenes and are spread across the entire building. Even though the layouts and appearances of the test buildings are different than the training buildings, our framework is able to adapt to the domain shift.

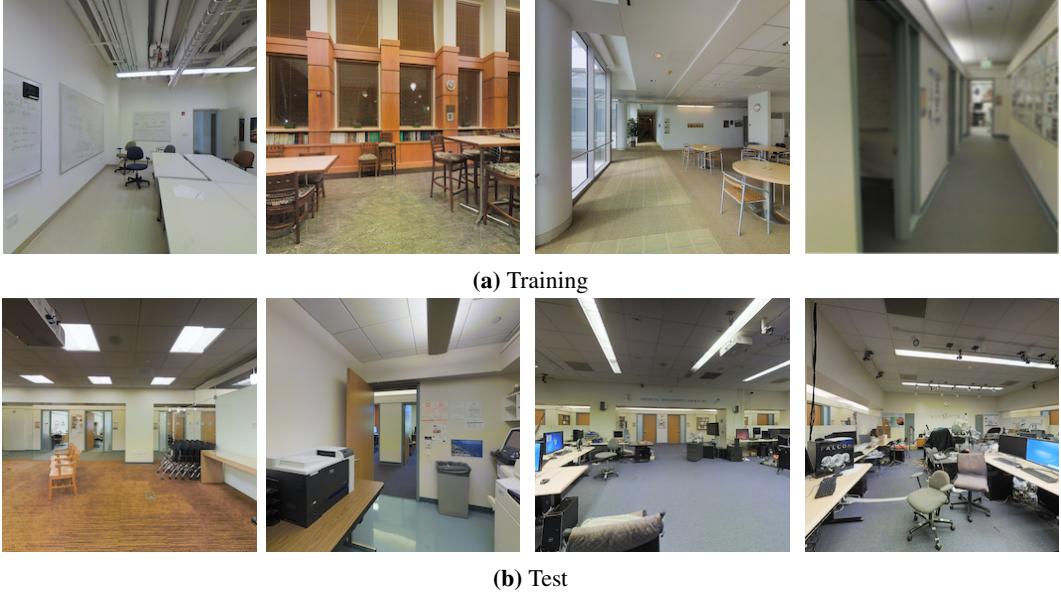


Figure 7. Some representative images of the buildings from which the training and the test data was collected. Even though the test environments are also office buildings, their layouts and appearances are different than the training buildings. However, our framework is still able to generalize to the domain shift.

In our experiments, we test the robot in two different buildings, none of which is in the training dataset (in fact, not even in the S3DIS dataset). We show some representative images of our experiment environments in Figure 8.



Figure 8. Some representative images of the buildings in which the experiments were conducted. None of these buildings were used for training/testing purposes in simulation.

8.5 Learning Navigation Affordances

In this section, we visualize some additional test cases and demonstrate how the proposed approach is indeed able to learn navigation cues for an efficient navigation in novel environments. In the first test case, the robot starts inside a conference room and is required to go to a target position that is in a different room. The problem setup is shown in Figure 9a. The blue dot represents the robot's start position and the black arrow represents the robot's heading. The center of the green area represents the target position.

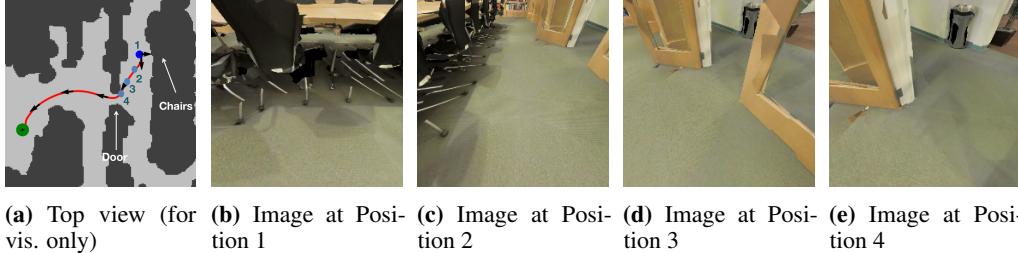


Figure 9. We visualize the trajectory as well as the observed RGB images (at the marked points) for LB-WayPtNav for a sample test task: the robot needs to go from the current room into another room. Our method is able to learn the navigation cue of exiting the room through the doorway to get to the goal location. Such cues enable the robot to do a more efficient and guided exploration in novel environments.

We mark some intermediate points on the robot trajectory in Figure 9a where it is required to predict a waypoint, and visualize the corresponding RGB images that it observed in Figures 9b-9e. The robot is initially facing some chairs (Fig. 9b), and based on the observed image, it predicts a waypoint to rotate in place to explore the environment. Upon rotation, the robot sees part of the door (Fig. 9c) and is able to learn that to go to the target position it needs to exit the current room through this door and predicts a waypoint to go towards the door (Fig. 9d). Next, it produces a waypoint to go out of the door (Fig. 9e), and eventually reaches the target position. We also show the full trajectory (the red plot) of the robot from the start position to the target position as well as mark the obstacles in context in Figure 9a.

As a second example, we consider the test case in Fig. 10a where the robot needs to go from Hallway 1 to Hallway 2. We show an intermediate point in the robot trajectory (Marked 1 in Fig. 10a) where the robot needs to turn right into Hallway 2 to reach the target. We also show the corresponding (stylized) field-of-view (projected on the ground) of the robot (the blue cone in Fig. 10a). The corresponding RGB image that the robot observed is shown in Figure 10b. Based on the observed image, the robot produces the cyan waypoint and the bold red desired trajectory. Even though the robot is not able to see its next desired position (it is blocked by the door in the field-of-view), it can see part of Hallway 2 and is able to reason that a hallway typically leads to another hallway. The robot uses this learned semantic prior to produce a waypoint leading into Hallway 2 and eventually reaches the target.

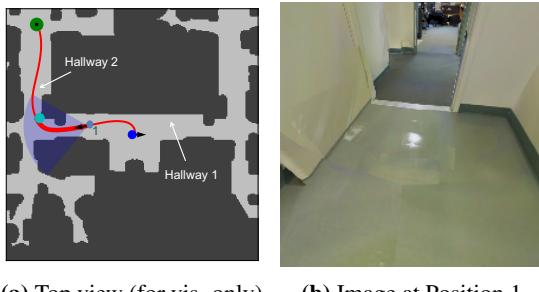


Figure 10. We visualize the trajectory as well as the observed RGB images (at the marked point) for our method (LB-WayPtNav) for a sample test task. In this task, the robot needs to go from one hallway (Hallway 1) into another (Hallway 2). Even though the entire Hallway 2 is not explicitly visible, the robot is able to reason that a hallway typically leads into another hallway, and uses this learned prior to efficiently reach its target position.

8.6 Hardware testbed

We use a TurtleBot 2 platform with a Yujin Kobuki robot base to serve as our autonomous vehicle during hardware experiments. The TurtleBot 2 is a low-cost, open source differential drive robot, which we equip with an Orbbec Astra RGB-D camera. However, for LB-WayPtNav and E2E learning experiments, we only used the RGB image. For geometric mapping and planning-based schemes, we additionally use the depth image. A snapshot of our testbed is shown in Figure 11.

The bulk of computation, including the deep network and the planning, runs on an onboard computer (Nvidia GTX 1060). The camera is attached to the onboard computer through a USB and supplies the RGB images. Given an image and the relative goal position, the onboard computer predicts the next waypoint for the robot, plans a spline trajectory to that waypoint, as well as computes the low-level control commands and the corresponding feedback controller. The desired speed and angular speed commands are sent to the Kobuki base, which then converts them to PWM signals to execute on the robot.

8.7 Imperfections in depth estimation

In Figure 12, we illustrate some examples of inaccurate depth estimations that we encounter during our experiments. The black pixels in the depth images correspond to the regions where the depth

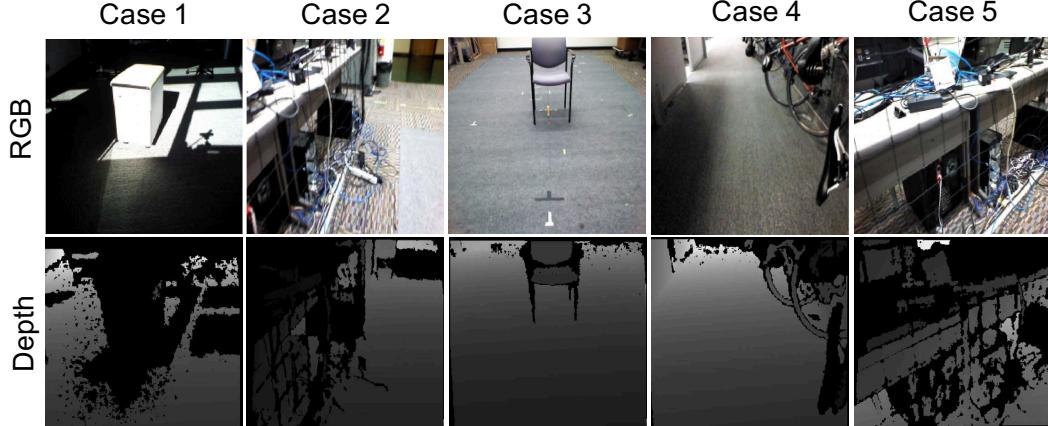


Figure 12. We visualize the RGB images captured by the robot and corresponding depth estimation. The black region corresponds to unknown depth. The depth estimation is inaccurate when the robot encounters shiny, transparent and thin objects, resulting in a significant decline in the performance of a mapping-based approach.

estimator fails to accurately estimate the depth. In Case 1, the depth estimation fails because of the presence of sunlight. In Case 2, thin wires and power strips on the floor are not recorded accurately. In Case 3, the depth estimation fails due to the shiny, thin legs of the chair. In Case 4, bike seats and tire frames are not estimated accurately. In Case 5, shiny, transparent monitors are not recognized by the depth estimator.



Figure 11. Our Turtlebot 2 hardware platform uses a Yujin Kobuki base, Gigabyte Aero Laptop, and Orbbec Astra camera.