# Learning to Navigate Sidewalks in Outdoor Environments

Maks Sorokin[1], Jie Tan[2], C. Karen Liu[3], Sehoon Ha[12]

*Abstract*— Outdoor navigation on sidewalks in urban environments is the key technology behind important human assistive applications, such as last-mile delivery or neighborhood patrol. This paper aims to develop a quadruped robot that follows a route plan generated by public map services, while remaining on sidewalks and avoiding collisions with obstacles and pedestrians. We devise a two-staged learning framework, which first trains a teacher agent in an abstract world with privileged ground-truth information, and then applies Behavior Cloning to teach the skills to a student agent who only has access to realistic sensors. The main research effort of this paper focuses on overcoming challenges when deploying the student policy on a quadruped robot in the real world. We propose methodologies for designing sensing modalities, network architectures, and training procedures to enable zero-shot policy transfer to unstructured and dynamic real outdoor environments. We evaluate our learning framework on a quadrupedal robot navigating sidewalks in the city of Atlanta, USA. Using the learned navigation policy and its onboard sensors, the robot is able to walk 3.2 kilometers with a limited number of human interventions.

Project webpage: **https://initmaks.com/navigation**

## I. INTRODUCTION

Outdoor navigation by foot in an urban environment is an essential life skill we acquired at a young age. Likewise, autonomous robotic workers that assist humans in an urban environment, such as last-mile delivery or neighborhood patrol, also need to learn to navigate sidewalks and avoid collisions. Similar to autonomous driving and indoor navigation, sidewalk navigation requires the robots to follow a route plan under noisy localization signals and egocentric sensors. However, outdoor sidewalk navigation faces more unstructured environments with a wide variety of pedestrians and obstacles and without any guiding lanes.

This paper aims to develop a learning framework that allows a quadrupedal robot to follow a route plan generated by public map services, such as Google or Apple Maps, while remaining on sidewalks and avoiding collisions. We take a two-staged learning approach similar to "learning by cheating" [1], where we first train a teacher agent in an abstract world with privileged information of bird-eye-view observations. Then clone the learned behavior to a student policy with realistic sensor configurations using Dataset Aggregation Method (DAGGER) [2].

However, directly applying "learning by cheating" [1] to outdoor navigation problem results in a poor policy that fails

[1] Georgia Institute of Technology, Atlanta, GA, 30308, USA
{maks,sehoonha}@gatech.edu
[2] Robotics at Google, Mountain View, CA, 94043, USA
jietan@google.com
[3] Stanford University, Stanford,CA, 94305, USA
karenliu@cs.stanford.edu

Fig. 1. AlienGo robot navigating various real-world sidewalks in outdoor.

immediately in the real world. Therefore, the main research contribution of this work is to methodologically ablate the sim-to-real gap and find pragmatic solutions to mitigate the primary sources of the gap, enabling [1] to be applied to the real robots.

Guided by careful analyses of the sim-to-real gap, we propose new methodologies for designing sensing modalities, network architectures, and training procedures. These inventions allow us to achieve the challenging goal of quadrupeds navigating outdoor sidewalks for a long distance with only a limited number of human interventions and without any curation or arrangement of the environment.

We evaluate the proposed framework on a real quadrupedal robot, AlienGo from Unitree, to walk the sidewalks in the city of Atlanta, USA. In our experiments, the robot navigated the total 3.2 kilometers of sidewalks using egocentric camera, LiDAR, and GPS with a few human interventions. Our natural testing environments include diverse static and dynamic objects typically seen in an urban space, such as sidewalks, trees, bridges, pedestrians, dogs, and bicycle riders. We also validate our framework by comparing different learning algorithms and sensor configurations in both simulated and real-world environments. Our technical contributions include:

- We develop a two-staged learning framework that leverages two different simulators, inspired by the "learning

by cheating" approach.

- We conduct a careful analysis of the *sim-to-real gap* and propose design choices and mitigation techniques.
- We demonstrate real-world outdoor sidewalk navigation on a quadrupedal robot, which spans 3.2 kilometers.

## II. RELATED WORK

### A. Visual Navigation

Researchers have studied vision-based robot navigation for decades using both hand-engineered and learning-based approaches. In 2007, Morales *et al.* [3] presented an autonomous robot system that can traverse previously mapped sidewalk trajectories in the Tsukuba Challenge. Kümmerle *et al.* [4] showed autonomous navigation of a 3 km route in the city center of Freiburg using the pre-built map data and a large array of laser sensors. While demonstrating impressive results, these systems require extensive manual engineering of various components, including sensor configuration, calibration, signal processing, and sensor fusion. As reported in their papers, it is also not straightforward to achieve complex behaviors, such as pedestrian avoidance.

On the other hand, deep learning enables us to learn effective navigation policies from a large amount of experience. The recent development of high-performance navigation simulators, such as Habitat [5] and iGibson [6], has enabled researchers to develop large scale visual navigation algorithms [7], [8], [9], [10], [11], [12], [13], [14] for indoor environments. Wijmans *et al.* [9] proposed an end-to-end vision-based reinforcement learning algorithm to train near-perfect agents that can navigate unseen indoor environments without access to the map by leveraging billions of simulation samples. Chaplot *et al.* [10] presented a hierarchical and modular approach for indoor floor construction & floor exploration that combines both learning and traditional algorithms. Francis *et al.* [11] showed an effective indoor navigation skill solely based on a single laser sensor.

While a lot of progress in learned indoor navigation has been made, only a handful of learning algorithms have been developed for outdoor navigation [15], [16], [17]. Müller *et al.* [15] trained a waypoint navigation policy in simulation and deployed it on a real toy-sized truck to drive on an empty vehicle road. Kahn *et al.* [17] proposed an approach for learning a sidewalk following policy purely in the real world by leveraging human-operator disengagements and overtakes. Instead, we demonstrate that it is possible to learn a policy to follow a designated set of waypoints on sidewalks without relying on any real-world data.

### B. Teacher-Student Framework

While end-to-end learning is a viable method, some challenging tasks require more structured learning algorithms, particularly when they involve large observation or state spaces that are hard to explore. Recently, Chen *et al.* [1] proposed the "*learning by cheating*" framework that can effectively learn complex policies by taking a two-staged approach. It first learns an expert policy that "cheats" by accessing privileged ground-truth information and transferring the learned behaviors with realistic sensor modalities. Lee *et al.* [18] demonstrated that this framework can learn a robust locomotion policy on challenging terrains.

We use the "*learning by cheating*" framework as a fundamental building block and demonstrate how we can further accelerate it by leveraging fast-speed simulation environments [12], [19]. While the original paper did not address the sim-to-real transfer, with our proposed techniques, we show that learning by cheating framework can be extended to real quadrupeds navigating in well-populated urban areas.

### C. Sim-to-Real Transfer

Overcoming the sim-to-real gap has been an active research topic in recent years [20], [21], [15], [22]. For instance, previous work [20], [21] investigated the Domain Randomization approach by applying random textures to the world surfaces to learn texture agnostic visual features. James *et al.* [22] proposed to train generative adversarial networks to reduce visual discrepancies between simulation and real-world images. Laskin *et al.* [23] proposed an image-based data augmentation technique to avoid overfitting, which applies random transforms to image observations.

Inspired by the work of Müller *et al.* [15], we also employ pretrained semantic networks for visual inference to help the sim-to-real transfer. However, we found that learned semantics networks with publicly available autonomous driving data sets show poor prediction on sidewalk navigation due to the perspective shift. Therefore, we retrain semantics with more sidewalk perspective images collected from virtual worlds [24]. We also train a policy using intermediate semantic features, as Shah *et al.* [25] suggested. Combining these techniques, we were able to bridge the visual sim-to-real gap for sidewalk navigation effectively.

## III. NAVIGATION POLICY

This section describes a learning algorithm to obtain a visual navigation policy for outdoor sidewalk navigation. Because of the complexity of the given problem, we adopt an approach of "learning by cheating" [1], [18], that decomposes the problem into learning of teacher and student policies (Fig. 2). First, we perform teacher training in an abstract world with salient information most relevant to the task of navigation. We use a birds-eye-view image as the privileged observation and efficiently learn a navigation policy. Once we obtain satisfactory performance, we train a student agent which only has access to realistic sensors (hence more limited sensing capability) in a high-fidelity simulator. We use DAGGER [2] train the student policy. This two-staged learning allows us to train a policy much more efficiently than learning it from scratch.

### A. Teacher Training in Abstract World

In the "learning by cheating" framework, teacher learning aims to efficiently obtain the ideal control policy by allowing it to access "privileged" information that cannot be obtained with robot's onboard sensors. In our scenario, privileged
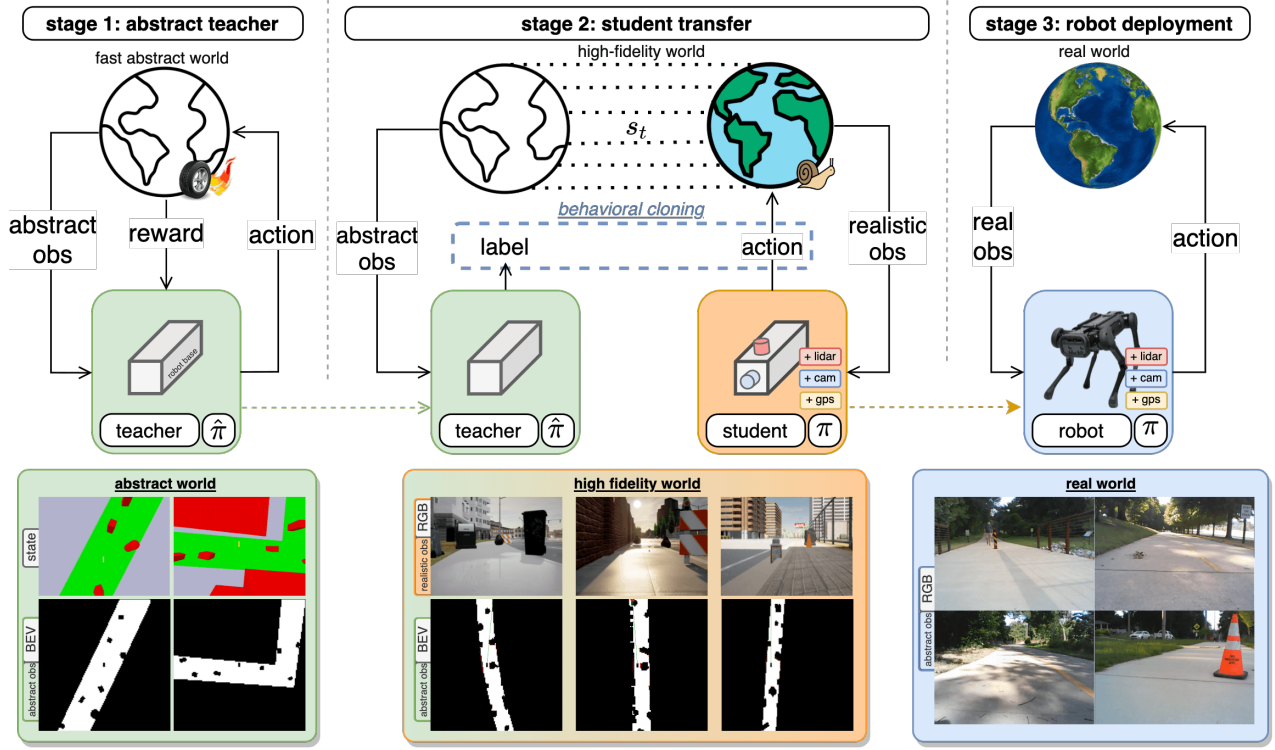
Fig. 2. Overview of the entire pipeline from abstract to real world deployment. A teacher is trained in an abstract world with privileged sensing information but simpler geometry and rendering thus enabling faster learning. A student is cloned in a high-fidelity simulation with realistic sensors, and finally deployed to the real world.

information is a birds-eye-view image that captures map layouts and nearby obstacles.

Because a high-fidelity simulator is computationally costly, we accelerate the learning by employing a simple abstract world for teacher training. This abstract world only contains essential information required for sidewalk navigation, such as walkable/non-walkable area, and static/dynamic objects (Fig. 2). We create this abstract world using Py-Bullet [26], which provides a fast and simple interface for accessing and rendering the abstract world observation modalities. In our experience, an abstract world can generate samples more than 10 times faster than a high-fidelity simulator, CARLA [27] because it has simplified geometry and avoids heavy rendering pipelines.

**Environment.** We generate the abstract world using the actual city data of the OpenStreetMap [28] with the help of OSM2World [29]. We parse the map layout information of the Helsinki area around $1km^2$ in size. To make the simulated empty street scenes closer to those in the real world, we randomly populate sidewalks with cylinder and cuboid-shaped objects. We also vary the sidewalk width from 2 to 5 meters for a richer training experience.

**Task.** The main objective of the teacher agent is to navigate the sidewalk to follow a list of waypoints on the map while avoiding collisions and respecting the sidewalk boundary. For each episode, we randomly select the start and goal sidewalk positions 10 to 15 meters apart while guaranteeing reachability. We design the reward function similar to other navigation papers [30]:

$$r = r_{success} + r_{termination} + r_{approach} + r_{life}.$$

$r_{success}$ is a sparse reward of 10.0 awarded if the agents gets close to the goal ($< 0.5$m). $r_{termination}$ is another sparse term of $-10.0$ if the agent collides with any obstacle, moves off the sidewalk, or fails to reach the goal within 150 simulation steps. $r_{approach}$ acts as a continuous incentive for making progress defined by the difference between the previous and the current distances to the goal. $r_{life}$ a small life penalty of $-0.01$ to avoid the agent from idling.

**Observation Space.** We define a *abstract* observation space (Fig. 2) with privileged information, which can be obtained from both the abstract and high-fidelity worlds. It consists of the following components:

1) **Bird-Eye-View Image** (BEV, privileged) is a 128x128 top-down binary image which covers a 18m by 18m region around the robot. Each pixel is one if the region is walkable while 0 for objects, walls, buildings, etc. We stack three previous BEVs to consider the history.

2) **Bird-Eye-View Lidar** (BLID) is a simulated LiDAR observation that measures distances to the nearest objects by casting 64 rays.

3) **Goal direction & distance** (GDD) is a relative goal location in the local polar coordinate frame.

**Action Space.** Each action is defined with two continuous values. The first action controls the speed (20 cm to $-10$ cm per step) in the forward/backward direction, while the second action controls the rotation ($-54°$ to $54°$ per step) around the yaw axis.

**Learning.** To train the teacher agent, we use a modification of the Soft Actor-Critic [31], [32]. We adjust the Actor and Critic networks to incorporate the BLID and GDD
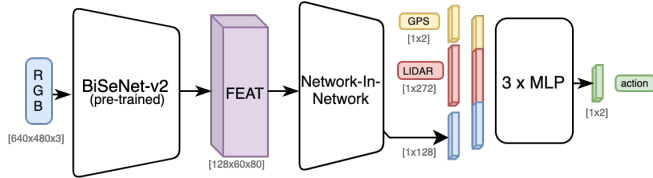
Fig. 3. A network architecture of the student policy.

| Agent | Train | Valid | Real World |
|---|---|---|---|
| **TEACHER** | 79.60% | 83.18% | – |
| **ULTIMATE** | 67.68% | 76.86% | – |
| **DEPTH** | 44.25% | 65.14% | – |
| **GT-SEM** | 67.02% | 73.14% | – |
| **INF-SEM** | 55.01% | 63.26% | 7.1% |
| **RGB** | **77.92%** | 73.16% | 25.0% |
| **INF-SEM-FEAT (ours)** | 69.71% | **77.90%** | **83.3%** |

TABLE I

<small>SUCCESS RATE COMPARISON OF DIFFERENT AGENT CONFIGURATIONS DURING TRAINING, VALIDATION, AND REAL WORLD TESTING.</small>

observations by stacking them with features extracted by a visual encoder.

We train the policy until saturation, which reaches the success rate of $83\%$ on the validation environments in the high fidelity world.

### B. Student Training in High fidelity world

Once we obtain an effective teacher policy, the next step is to learn a student policy in a high-fidelity simulator with realistic observations that the robot's onboard sensors can provide. We train a student agent by cloning teacher's behavior using Dataset Aggregation Method (DAGGER) [2]. **Environment.** We train a student policy in CARLA [27] (Fig. 2), a high fidelity simulator built on the Unreal Engine 4. The simulator includes multiple town blocks with various weather conditions, which in total is $0.6$ km$^2$ big. Similar to the abstract world, we densely populate random objects on sidewalks to increase the difficulty. In total, we use 7 towns and 14 weather conditions for training and validation, and test the best policy on the real robot.

**Observation Space.** Our AlienGo robot is equipped with the following onboard sensors: an RGB camera, a depth camera, a 1-D LiDAR, a GPS (from a cell phone), and a T265 camera for visual odometry. Additionally, we can generate semantic segmentation of the scene by processing RGB images with pretrained networks. Designing an observation space that enables sim-to-real transfer is a key contribution of this work and will be detailed in Section IV. Here we list the observation modalities we use:

1) **Real Feature** (RFEAT) is a 128x60x80 feature tensor extracted from RGB image using pre-trained semantic network.
2) **Real LiDAR** (RLID) is a 272 LiDAR information reading, representing the distance to the nearest object captured by casting a laser in all directions. The maximum sensing distance is capped at 6 meters.
3) **Real Goal direction & distance** (RGDD) is the same relative goal location in the local polar coordinate frame, which is computed based on GPS signals.

**Behavior Cloning with DAGGER.** We clone the teacher's behavior to the student policy using DAGGER. DAGGER allows the student agent to learn proper actions in a supervised learning fashion in the presence of the expert. It generates a rollout using the current student policy, collects the corresponding *abstract* and *realistic* observations, and updates the parameters by minimizing the L1 loss between the student and teacher actions. Note that in the high-fidelity

simulator we can generate both *realistic* and *abstract* observations, where the latter is used by the teacher for generating labels (Fig. 2). We bootstrap the learning by prefilling the replay buffer with $120k$ successful teacher experience pairs before the student training begins. The student architecture (Fig. 3) consists of 3 components: pre-trained BiSeNet [33] without the last layer for feature extraction, Network-In-Network [34], and 3-layer MLP for feature processing.

In our experience, our two-staged learning framework takes approximately 30 hours to converge, while learning from scratch takes more than 300 hours to converge.

## IV. SIM-TO-REAL TRANSFER

After learning in simulation, the next step is to deploy the policy to the real quadruped. While the *learning by cheating* framework offers sample efficient learning, it can suffer from severe performance degradation during both behavior cloning and real-world transfer. To cross the *sim-to-real* gap, we need to carefully design sensing modalities of the policy and data augmentation method to train the semantic model.

### A. Sensing Modalities

In simulated environments, adding more sensing modalities often translates to simpler learning problems [35], [36]. However, in the real world, we need to cautiously select sensors because each sensor comes with its own *sim-to-reap* gap. Therefore, we evaluate each sensor based on three criteria: (1) the usefulness of the information it encapsulates, (2) the additional difficulty in learning it adds, and (3) the sensitivity to the sim-to-real gap it induces.

To this end, we conduct an ablation study (Table I) of different agent configurations. Our simulation experiments help us to measure the usefulness (criteria 1) and learning difficulty (criteria 2), while real-world experiments evaluate sim-to-real transferability (criteria 3). For more details of real-world experiments, please refer to Figure 4 and Section V-A. Based on the experimental results, we select three sensor modalities: inferred semantic features (RFEAT), lidar (RLID), and goal direction and distance (RGDD). We will describe the selection process here.

**Sanity Check.** First, we conduct a sanity check to test whether the teacher's policy can be successfully transferred to the student agent with egocentric observations. For this purpose, we train an ideal *ULTIMATE* agent that has access

to all the sensors, raw RGB images, depth images, ground-truth semantic images, lidar, and localization. It achieves the success rate of 76.86% that is only 6% lower than the teacher (Table I). Therefore, we conclude it a successful sanity check.

**Visual Sensors.** Visual sensors, such as RGB, depth images, or semantics segmentations, are useful information to recognize surroundings. Table I indicates that all three agent with individual sensors, raw RGB images (*RGB*), depth images (*DEPTH*), ground-truth semantic images (*GT-SEM*), show promising performance of 73.16%, 65.14% and 73.14%, respectively. However, ground-truth semantic is not available at deployment, and the performance drops to 63.26% when we infer it using a pre-trained semantic segmentation model (*INF-SEM*). Then we investigate the third criteria, the sim-to-real transferability. At this point, we stop investigating the *DEPTH* agent because our infrared-based depth camera works poorly under direct sunlight in outdoor environments. Unfortunately, both *INF-SEM* and *RGB* agents show poor success rates of 7.1% and 25% in the real world. While the RGB agent seems a bit more promising, we are not able to improve its performance in the real world, after applying commonly-used sim-to-real techniques, such as domain randomization [23].

The significant performance degradation of *INF-SEM* in the real world is due to the poor generalization of the semantic segmentation model. Our insight is that the last layer of the model discards important features to classify each pixel into a small number of classes. Therefore, we suspect that the features before the output layer may contain more useful information than the final labels.

For this reason, we trained a new agent, *INF-SEM-FEAT*, which takes the *feature layer* (the last layer before prediction) of the semantic network as the input. This approach is common in transfer learning in computer vision [25], but under-explored in visual navigation [9]. This agent achieves satisfactory success rates in both simulation and the real world, 77.90% and 83.3%, respectively.

**LiDAR.** We select a LiDAR as an additional sensor because it contains useful distance information for collision avoidance and suffers less from noise in the real world. Therefore, we train all the agents with an additional LiDAR sensor.

**Localization.** We test two common approaches for outdoor localization: visual odometry and GPS. In our experience, visual odometry, such as the Intel Realsense T265 sensor, is vulnerable to the accumulation of errors, and thus is not desirable for long-range navigation tasks. Therefore, we decide to use a GPS sensor on a Pixel5 phone, despite its high latency and noise levels. We find that the phone's GPS works poorly when it is mounted on our legged robot. This is probably due to the high frequency vibrations of the robot's body, caused by periodic foot impacts during walking. Instead of using a more advanced GPS or developing a damping mechanism, which are orthogonal to the goal of this paper, we simply carry the phone with a human operator who walks closely with the robot.

## B. Data Augmentation

The performance of semantic agents, *INF-SEM* and *INF-SEM-FEAT* highly depends on pre-trained semantic networks. However, when we tested publicly available segmentation networks, such as [33], [37], their performance is not satisfactory because they are mainly trained with autonomous-driving datasets and suffer from perspective shifts [38] from road to sidewalk. Therefore, we take the BiSeNetV2 [33] implementation by CoinCheung [39] and slightly adapt it for our task.

To expand training set with additional sidewalk images, we create an augmented dataset that contains 17535 sidewalk perspective and 9704 road images synthesized using GTA [24] and 6000 road and 2400 sidewalk images synthesized using CARLA. We combine these newly generated synthesized images with the existing 2975 real images of Cityscapes [40] for training.

However, this data augmentation induces new sidewalk perspectives to the existing dataset, which previously had only road perspectives, and makes the learning of BiSeNet significantly more difficult. This is because the network tends to leverage the location of the surface as a strong prior. In other words, the original BiSeNet tends to predict the surface under the camera as a vehicle road. It is true for vehicle perspectives, but not true for sidewalk perspectives.

We address this issue by solving a slightly different problem of segmenting the region a robot currently occupies. For example, if the robot is standing on the road, we segment the entire road in the field of view. Likewise, if the robot is currently walking on the sidewalk, the segmentation region will be the sidewalk. Once trained, we evaluated the model on 35 manually annotated images from the real robot's RGB camera and observed clear improvements.

## V. EXPERIMENTS

Three key requirements of navigating sidewalks are 1) staying on the sidewalk, 2) avoiding collisions with obstacles and pedestrians and 3) reaching the goal. To validate that our approach satisfies all three requirements, we perform two types of real-world experiments. In Section V-A, we evaluate the ability of our method (*INF-SEM-FEAT*) and two baselines (*RGB* and *INF-SEM*) on obstacle avoidance and sidewalk following. In Section V-B, we evaluate our best performing agent, *INF-SEM-FEAT*, on long-range sidewalk navigation tasks, where the robot gets stress-tested in a wide variety of real-world scenarios.

## A. Obstacle Avoidance

**Experiment Setup.** We select a random sidewalk section of 15 meters in length at the nearby park as a testing ground. We populate a variety of static obstacles represented that can be commonly found in a household (Figure 4 Left). Neither the sidewalk nor the objects were previously seen by the policy during training.

The agent is instructed to walk along the sidewalk towards the other end. We disable GPS and the robot only relies on visual odometry for this experiment as the agent does not need
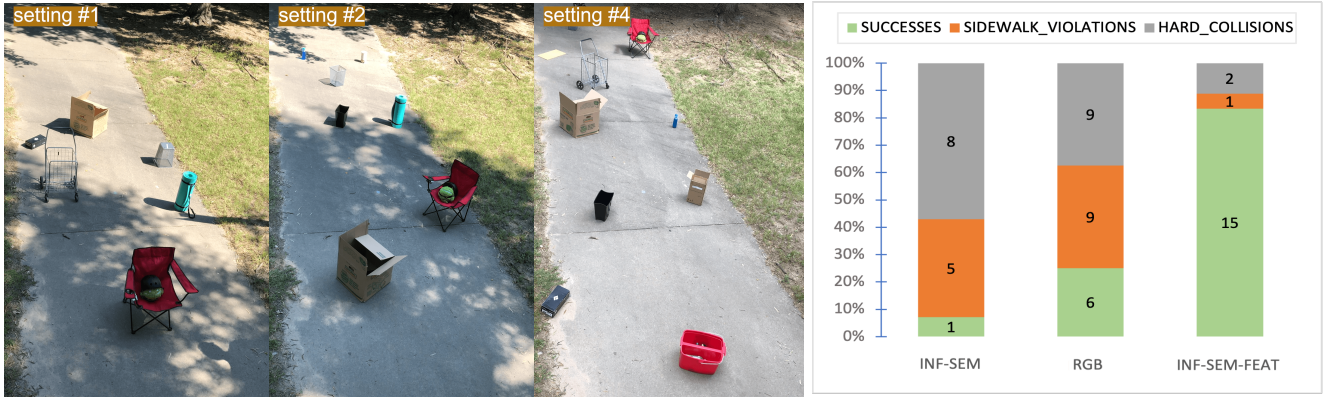
Fig. 4. Comparison of the collision avoidance performance of three agents: *RGB*, *INF-SEM*, and *INF-SEM-FEAT*. **Left:** Two out of four real world course configurations **Right:** Success rates with the breakdown of failure types: *sidewalk violations* and *hard collisions*. Our agent, *INF-SEM-FEAT*, shows a much higher success rate compared to other two agents.

to walk far and is only guided by directional information. The robot is placed at the starting location and runs until it triggers one of the termination conditions, judged by a human supervisor. Three possible outcomes are (1) *success* when the agent successfully passes the obstacle course without any violation (2) *collision* when the agent collides with obstacles and (3) *sidewalk violation* when the agent steps out of the sidewalk. We ignore minor collisions that do not noticeably affect the robot's trajectory or obstacles' positions.

**Results.** We test all three agents multiple times on four different sidewalk settings (Figure 4 Left). We report *success rate* of each agent in Figure 4. Right along with the breakdown of failure reasons: *hard collisions* or *sidewalk violations*. In agreement with our finding in simulation, the *INF-SEM-FEAT* achieves the highest success rate (83.3%) when deployed in the real world. It only collides with obstacles only twice and violates the sidewalk once. the *INF-SEM* and *RGB* agents show much lower success rates, 7.1% and 25.0%, respectively. We also observed that the *RGB* agent shows a slightly lower *hard collisions rate* than the *INF-SEM* agent.

Both *RGB* and *INF-SEM* agents are more sensitive to variations in sunlight, as they are both attracted to well-lit regions resulting in sidewalk violation. In our experience, *INF-SEM-FEAT* is more robust to variances in lighting conditions or textures. We conjecture that this is due to richer information of the feature layer compared to the thresholded semantic predictions or raw RGB pixels.

For minor collisions that did not trigger termination, we checked the recorded sensor data, and found that these collisions were often caused by sensor limitations, e.g. the object is behind the RGB camera or too small/thin for the LiDAR to detect. These collisions can be potentially mitigated by longer observation history or additional sensing.

### B. Long-distance Sidewalk Navigation

**Experiment Setup.** To comprehensively evaluate our method, we deploy the learned navigation policy and the robot on a wide variety of real-world sidewalks and trails. We define six different routes in three locations around the city



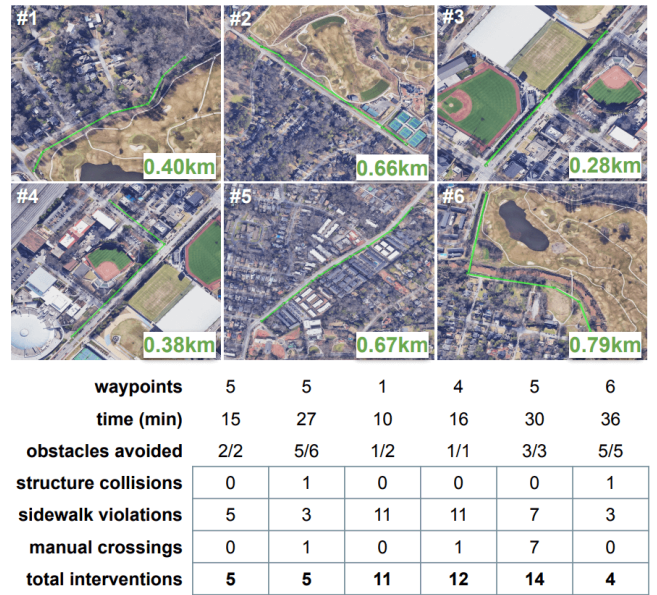| | #1 | #2 | #3 | #4 | #5 | #6 |
|---|---|---|---|---|---|---|
| waypoints | 5 | 5 | 1 | 4 | 5 | 6 |
| time (min) | 15 | 27 | 10 | 16 | 30 | 36 |
| obstacles avoided | 2/2 | 5/6 | 1/2 | 1/1 | 3/3 | 5/5 |
| structure collisions | 0 | 1 | 0 | 0 | 0 | 1 |
| sidewalk violations | 5 | 3 | 11 | 11 | 7 | 3 |
| manual crossings | 0 | 1 | 0 | 1 | 7 | 0 |
| total interventions | **5** | **5** | **11** | **12** | **14** | **4** |

Fig. 5. Our robot is able to autonomously navigate six routes using its onboard sensors, with limited human supervision. We report a few statistics about the routes and the breakdown of interventions.

of Atlanta, USA, with total length of 3.2 kilometers (Fig. 5). We select a set of GPS waypoints along each route. When the robot reaches within the 2m radius of a given waypoint, it proceeds to the next one.

**Interventions.** Human supervisors constantly monitor the robot and take over control in the following three cases:

1) Safety: The robot walks too close to people around, or the robot gets into a dangerous situation for itself (e.g. imminent collisions).
2) Sidewalk violation: The robot walks out of the side-walk boundary and does not automatically recover.
3) Road crossing: The robot needs help crossing the road, as it is not yet trained to do so.

**Results & Discussion.** We summarize the numbers of how the policy performed on Fig. 5. Over 3.2 kilometer walk, the robot avoided 17 out of 19 obstacles. The longest non-

interrupted section was 320 meters, while the mean is from 31m to 198m, which varies to different courses.

*Obstacle avoidance:* The robot was exposed to naturally presented obstacles, such as poles, cones, and trash bins. Additionally, the human supervisors sometimes intentionally walked in front of the robot to pretend as pedestrians. The agent managed to avoid 17 out of 19 obstacles. The two collisions were against a trash bin and a narrow poll, which were minor and did not require human interventions. In both cases, the robot attempted to avoid the obstacle but collided with its side while passing around.

*Sidewalk following:* The agent generally possesses a good ability to stay on sidewalks. Although the robot sometimes failed at localization, it could still remain within the walkable regions. This robust behavior indicates that the agent is not only directed to the goal location but prioritizes staying in the sidewalk boundary. However, the agent sometimes still walked outside of sidewalk boundaries when it encountered visually challenging scenarios, such as bright regions or transitions to shadows.

Thanks to the rich features in the semantic model, the agent was able to navigate on different sidewalk surfaces, including pavement, tiled, and many others available in the simulation. However, we observed that certain types of terrains could cause confusions to the robot. For instance, we found that the agent was more prone to violate the sidewalk boundary on two particular routes. We analyzed these failure cases by reconstructing semantic images from the recorded features. It turned out that the agent misclassified grassy sidewalks as non-walkable areas.

*Localization Errors:* In general, GPS performance was stable throughout the experiments. However, GPS occasionally experienced significant delays or noises, which caused errors in the localization, waypoint switching logic and ultimately led to sidewalk violations.

*Driveways:* Another main reason for manual interventions was the existence of driveways, a slope on the sidewalk to give car access to buildings. Note that driveways were not included in the training data. The agent often confused the semantics of the sidewalk and driveway due to the lack of data, and then proceeded in the direction of the driveway. Please refer to the supplemental video for examples.

**Q-function Analysis.** Along with an actor, we also learn a Q function. When visualizing the Q function together with semantic maps, we find that the changes of Q-values are interpretable and intuitive. In Fig. 6, we illustrate a few interesting situations, such as sidewalk violation, encountering a pedestrian, and unclear vision due to sun glare. All these scenarios result in decreasing Q-values. Once the situations improve, the Q-values recover back to the normal range.

## VI. CONCLUSION

We developed a quadrupedal robot that follows the route plan generated by the public map service, while remaining on sidewalks and avoiding collisions with obstacles and pedestrians using its onboard sensors. The main research effort focused on deploying a policy trained using a two-staged learning framework in simulation, to an unstructured and dynamic real outdoor environment. We discussed the important decisions on sensing modalities and presented new training procedures to overcome the sim-to-real gap. We evaluated our learning on a quadrupedal robot, which was able to walk 3.2 kilometers on sidewalks in the city of Atlanta, USA, with the learned navigation policy, onboard sensors, and a limited number of human interventions.

In the future, we plan to address a few limitations. First, the algorithm determines the waypoint transition moment based on an overly simplified rule-based logic that is vulnerable to localization errors. One potential approach is to leverage perceived semantics to determine transition moments, such as "turning at the corner in front of the stop sign". Second, the agent is subject to minor collisions when it passes the object due to the combination of limited sensing capability and lack of memory. In future work, we plan to address this issue by investigating stateful models, such as recurrent neural networks. Finally, we want teach the robot dog to cross roads automatically, while it is now manually controlled by a human operator. This would require us to have traffic lights and crosswalks in simulation, as well as active control of vision [41].

## REFERENCES

[1] D. Chen, B. Zhou, V. Koltun, and P. Krähenbühl, "Learning by cheating," in *Conference on Robot Learning*. PMLR, 2020.

[2] S. Ross, G. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011.

[3] Y. Morales, A. Carballo, E. Takeuchi, A. Aburadani, and T. Tsubouchi, "Autonomous robot navigation in outdoor cluttered pedestrian walkways," *Journal of Field Robotics*, vol. 26, no. 8, pp. 609–635, 2009.

[4] R. Kümmerle, M. Ruhnke, B. Steder, C. Stachniss, and W. Burgard, "Autonomous robot navigation in highly populated pedestrian zones," *Journal of Field Robotics*, vol. 32, no. 4, pp. 565–589, 2015.

[5] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik, D. Parikh, and D. Batra, "Habitat: A Platform for Embodied AI Research," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.

[6] B. Shen, F. Xia, C. Li, R. Martín-Martín, L. Fan, G. Wang, C. Pérez-D'Arpino, S. Buch, S. Srivastava, L. P. Tchapmi, M. E. Tchapmi, K. Vainio, J. Wong, L. Fei-Fei, and S. Savarese, "igibson 1.0: a simulation environment for interactive tasks in large realistic scenes," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2021, p. accepted.

[7] S. Gupta, J. Davidson, S. Levine, R. Sukthankar, and J. Malik, "Cognitive mapping and planning for visual navigation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2616–2625.

[8] K. Fang, A. Toshev, L. Fei-Fei, and S. Savarese, "Scene memory transformer for embodied agents in long-horizon tasks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 538–547.

[9] E. Wijmans, A. Kadian, A. Morcos, S. Lee, I. Essa, D. Parikh, M. Savva, and D. Batra, "Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames," *arXiv*, pp. arXiv–1911, 2019.
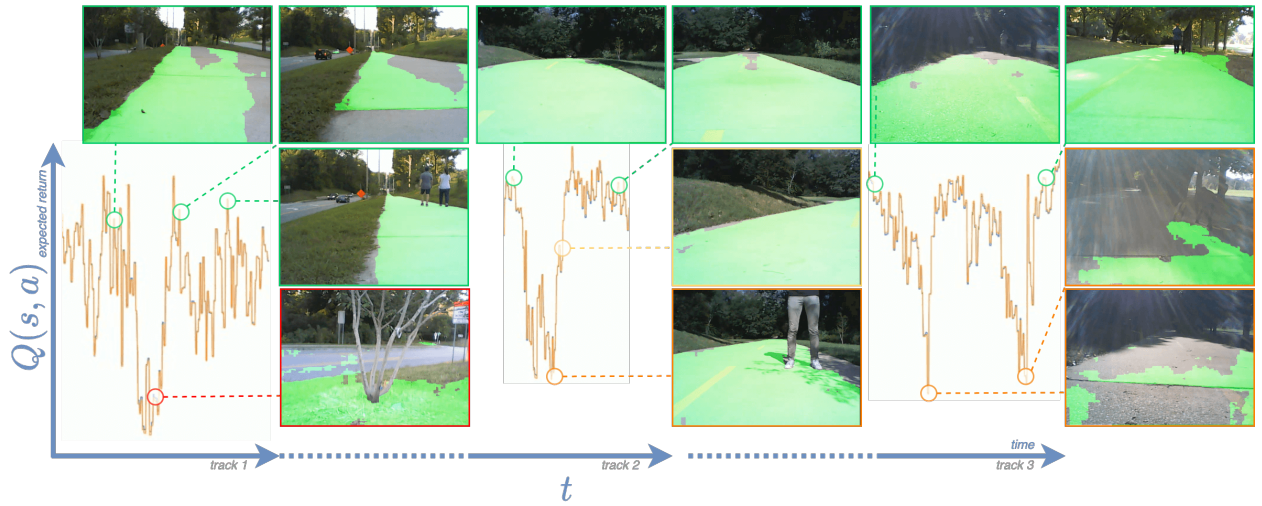
Fig. 6.    Q-function value predictions during real-world testing. Q-function values drop when an agent detects anomalies, such as sidewalk violation (**Left**), a human obstacle (**Middle**), and sun glare (**Right**), and recover back when the situation is resolved.

[10] D. S. Chaplot, D. Gandhi, S. Gupta, A. Gupta, and R. Salakhutdinov, "Learning to explore using active neural slam," *arXiv preprint arXiv:2004.05155*, 2020.

[11] A. Francis, A. Faust, H.-T. L. Chiang, J. Hsu, J. C. Kew, M. Fiser, and T.-W. E. Lee, "Long-range indoor navigation with prm-rl," *IEEE Transactions on Robotics*, vol. 36, no. 4, pp. 1115–1134, 2020.

[12] A. Petrenko, Z. Huang, T. Kumar, G. Sukhatme, and V. Koltun, "Sample factory: Egocentric 3d control from pixels at 100000 fps with asynchronous reinforcement learning," in *International Conference on Machine Learning*.   PMLR, 2020, pp. 7652–7662.

[13] D. Hoeller, L. Wellhausen, F. Farshidian, and M. Hutter, "Learning a state representation and navigation in cluttered and dynamic environments," *IEEE Robotics and Automation Letters*, 2021.

[14] H. Karnan, G. Warnell, X. Xiao, and P. Stone, "Voila: Visual-observation-only imitation learning for autonomous navigation," *arXiv preprint arXiv:2105.09371*, 2021.

[15] M. Müller, A. Dosovitskiy, B. Ghanem, and V. Koltun, "Driving policy transfer via modularity and abstraction," *arXiv preprint arXiv:1804.09364*, 2018.

[16] G. Kahn, P. Abbeel, and S. Levine, "Badgr: An autonomous self-supervised learning-based navigation system," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1312–1319, 2021.

[17] ——, "Land: Learning to navigate from disengagements," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1872–1879, 2021.

[18] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning quadrupedal locomotion over challenging terrain," *Science robotics*, vol. 5, no. 47, 2020.

[19] A. Petrenko, E. Wijmans, B. Shacklett, and V. Koltun, "Megaverse: Simulating embodied agents at one million experiences per second," in *ICML*, 2021.

[20] F. Sadeghi and S. Levine, "Cad2rl: Real single-image flight without a single real image," *arXiv preprint arXiv:1611.04201*, 2016.

[21] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*.   IEEE, 2017.

[22] S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell, and K. Bousmalis, "Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.

[23] M. Laskin, K. Lee, A. Stooke, L. Pinto, P. Abbeel, and A. Srinivas, "Reinforcement learning with augmented data," *arXiv preprint arXiv:2004.14990*, 2020.

[24] S. R. Richter, Z. Hayder, and V. Koltun, "Playing for benchmarks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2213–2222.

[25] R. Shah and V. Kumar, "Rrl: Resnet as representation for reinforcement learning," in *Self-Supervision for Reinforcement Learning Workshop-ICLR 2021*, 2021.

[26] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," http://pybullet.org, 2016–2021.

[27] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "Carla: An open urban driving simulator," in *Conference on robot learning*.   PMLR, 2017, pp. 1–16.

[28] OpenStreetMap contributors, "Planet dump retrieved from https://planet.osm.org ," https://www.openstreetmap.org, 2017.

[29] "Osm2world." [Online]. Available: https://github.com/tordanik/OSM2World

[30] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik *et al.*, "Habitat: A platform for embodied ai research," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9339–9347.

[31] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel *et al.*, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905*, 2018.

[32] D. Yarats, A. Zhang, I. Kostrikov, B. Amos, J. Pineau, and R. Fergus, "Improving sample efficiency in model-free reinforcement learning from images," *arXiv preprint arXiv:1910.01741*, 2019.

[33] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang, "Bisenet: Bilateral segmentation network for real-time semantic segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 325–341.

[34] M. Lin, Q. Chen, and S. Yan, "Network in network," *arXiv preprint arXiv:1312.4400*, 2013.

[35] B. Zhou, P. Krähenbühl, and V. Koltun, "Does computer vision matter for action?" *arXiv preprint arXiv:1905.12887*, 2019.

[36] J. T. Kim and S. Ha, "Observation space matters: Benchmark and optimization algorithm," *IEEE International Conference on Robotics and Automation*, 2021.

[37] R. Mohan and A. Valada, "Efficientps: Efficient panoptic segmentation," *International Journal of Computer Vision (IJCV)*, 2021.

[38] R. Shetty, B. Schiele, and M. Fritz, "Not using the car to see the sidewalk–quantifying and controlling the effects of context in classification and segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.

[39] "Bisenet by coincheung." [Online]. Available: https://github.com/CoinCheung/BiSeNet

[40] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.

[41] M. Sorokin, W. Yu, S. Ha, and C. K. Liu, "Learning human search behavior from egocentric visual inputs," in *Computer Graphics Forum*, vol. 40, no. 2.   Wiley Online Library, 2021, pp. 389–398.