

# On Embodied Visual Navigation in Real Environments Through Habitat\*

Marco Rosano<sup>1,3</sup>, Antonino Furnari<sup>1</sup>,  
Luigi Gulino<sup>3</sup>, Giovanni Maria Farinella<sup>1,2</sup>

<sup>1</sup>FPV@IPLAB - Department of Mathematics and Computer  
Science, University of Catania, Italy

<sup>2</sup>Cognitive Robotics and Social Sensing Laboratory,  
ICAR-CNR, Palermo, Italy,

<sup>3</sup>OrangeDev s.r.l., Firenze, Italy

marco.rosano@unict.it, furnari@dmi.unict.it,  
luigi.gulino@orangedev.it, gfarinella@dmi.unict.it

## Abstract

Visual navigation models based on deep learning can learn effective policies when trained on large amounts of visual observations through reinforcement learning. Unfortunately, collecting the required experience in the real world requires the deployment of a robotic platform, which is expensive and time-consuming. To deal with this limitation, several simulation platforms have been proposed in order to train visual navigation policies on virtual environments efficiently. Despite the advantages they offer, simulators present a limited realism in terms of appearance and physical dynamics, leading to navigation policies that do not generalize in the real world. In this paper, we propose a tool based on the Habitat simulator which exploits real world images of the environment, together with sensor and actuator noise models, to produce more realistic navigation episodes. We perform a range of experiments to assess the ability of such policies to generalize using virtual and real-world images, as well as observations transformed with unsupervised domain adaptation approaches. We also assess the impact of sensor and actuation noise on the navigation performance and investigate whether it allows to learn more robust navigation policies. We show that our tool can effectively help to train and evaluate navigation policies on real-world observations without running navigation episodes in the real world.

## 1 Introduction

The autonomous visual navigation problem has been studied for years by the research community and has recently attracted even more interest given its po-

---

\*Published in International Conference on Pattern Recognition (ICPR), 2020.

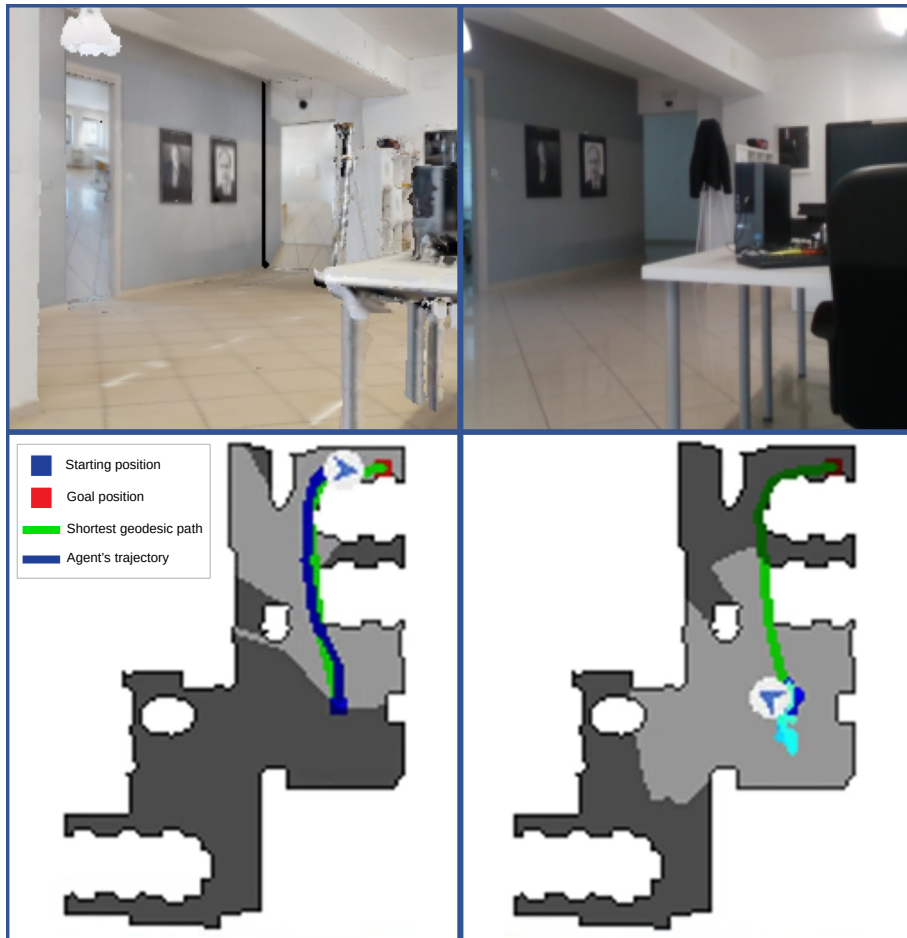


Figure 1: Navigation examples performed by an agent trained in a virtual environment when tested on (left) virtual observations and (right) real-world observations. The top row shows examples of observations seen during test. The policy in the virtual environment does not generalize to real-world observations.

tential range of applications in real contexts [1]. Recently, Deep Learning (DL) approaches have shown that it is possible to learn a navigation policy in an end-to-end fashion directly from visual observations collected by the agent while interacting with the environment [2, 3]. This approach outperforms past navigation paradigms and can learn a more abstract representation of the environment which allows to easily transfer the acquired knowledge to unseen contexts [4, 5]. One of the major drawbacks of DL-based navigation approaches is that they require to run several navigation episodes in order to learn effective navigation policies. Unfortunately, collecting the required experience in the real-world is difficult due to the maintenance cost required by a robotic platform acting in a trial-and-error setup, the variety of the scenes to be explored to ensure regularization, as well as the time needed to perform all the training episodes. A popular approach to overcome these limitations consists in collecting all the

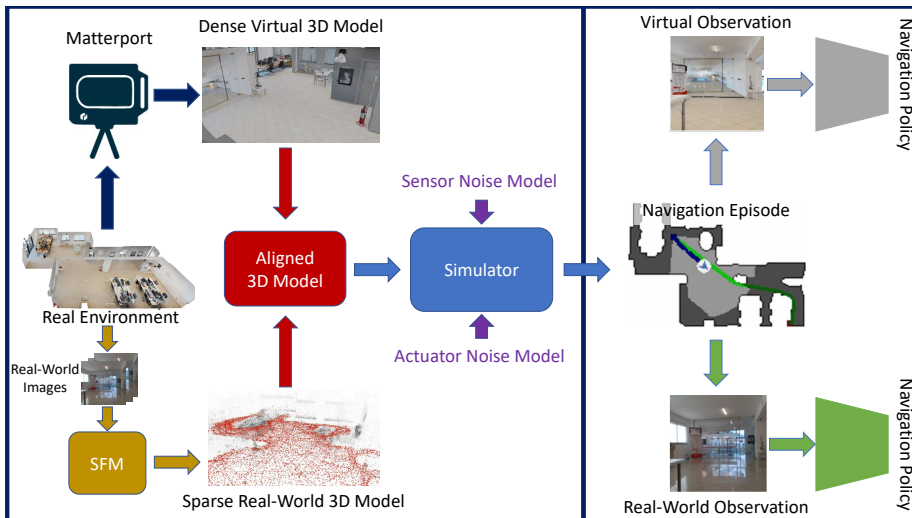


Figure 2: Our method exploits two 3D models of the same environment to generate realistic navigation episodes comprising both virtual and real-world observations, which can be used to evaluate or train navigation policies. The model can also account for sensor and actuator noise.

required experience in simulated environments which replicate the appearance of real scenes. Following this idea, several simulators for training embodied agents have been proposed in previous works [4, 6, 7]. These simulators allow to import 3D models of realistic environments previously acquired with a 3D scanner, resulting in a fast increase of the available training environments.

Even if training navigation agents with these tools has been shown to be effective, transferring the learned policies to the real world (e.g. to a real robotic platform) is still a challenging and open problem. This is mainly due to a domain shift caused by a number of factors such as the difference in the appearance between real-world visual observations and the ones provided by the simulator, the absence of real-world noise in the simulation (i.e., actuator noise and sensor noise), as well as the simplified physical interactions with the elements of the scene (e.g. no collisions, no friction, no uneven ground, etc.). Indeed, a navigation policy learned in a simulated environment performs very poorly when transferred to a real robotic agent who has to explore the environment [8]. Figure 1 shows two navigation examples obtained by an agent trained with virtual observations. When the agent is tested on visual observations coming from the simulator (Fig. 1 left), the navigation is close to perfect, whereas when the same agent is tested with real-world images (Fig. 1 right), the navigation fails. In this paper, we investigate how a navigation policy learned in a simulated world transfers to a corresponding real environment. Previous works have tried to measure and compensate for this domain shift by fine-tuning models learned in simulation using few real observations sampled on a sparse grid [2], or by running parallel experiments in simulation and in reality [8, 9]. While these approaches are viable, they considered simplified virtual environments [2] or require access to a real robotic platform during training and test [8, 9]. We argue that this prevents the exploitation of the convenient abstraction provided by a

simulator. Hence, we propose to perform the analysis by relying on a simulator to sample real-world observations, which allows to take advantage of the fast and inexpensive training and evaluation provided by modern simulators. The core of our approach is a tool based on the popular Habitat simulator [4] which allows to generate navigation episodes comprising virtual and real-world observations. The scheme of the approach is illustrated in Figure 2. Given a real environment, we construct two 3D models: 1) a “virtual” model using the Matterport 3D scanner<sup>1</sup>, and 2) a “real-world” model running a Structure from Motion algorithm (SfM) [10] on real images of the environment. The “virtual” model is an accurate representation of the world with limited photo-realism, whereas the “real-world” model is a sparse collection of real images attached with their related camera pose. The two models are then geometrically registered so that their coordinate systems match to form an “aligned 3D model”. The proposed tool hence allows to generate paired virtual and real-world navigation episodes by sampling visual observations from the aligned 3D model. Sensor and actuator noise models can be provided to enable the simulator to generate more realistic navigation episodes. The developed tool allows to 1) assess the performance of a navigation model which has been trained with virtual observations, when tested on real-world visual inputs; 2) measure the impact of sensors noise, actuator noise and appearance gap on a navigation policy learned with virtual observations when tested on real-world visual inputs; 3) adapt a navigation model learned with virtual observations to the real domain by fine-tuning it with real-world observations; 4) train realistic navigation policies introducing simulated actuator and sensor noise akin to the one which can be observed in real robotic platforms.

We perform experiments on an office environment of  $150m^2$ . Our analysis points out that navigation policies trained with virtual observations perform poorly when tested in real-world scenes, noisy sensors and noisy actuators have a significant impact on the navigation performance, models trained in a noisy environment can learn navigation policies which are more robust to noise. We also benchmark different approaches to improve the performance of navigation agents when deployed to the real world. We found that both unsupervised domain adaptation such as CycleGAN [11] and supervised adaptation such as fine-tuning, allow to adapt the navigation policy to the real-world scenario, with the latter class of approaches obtaining better performance. Results suggest that, while the approaches to domain adaptation for visual navigation are promising, much work is still to be done to achieve results which are exploitable in the real world. We hope that the proposed tool can foster research in this direction.

In sum, the contributions of this work are as follows:

1. we propose a tool based on Habitat [4] to train and evaluate entirely in simulation visual navigation policies on real observations and with realistic sensor and actuator noise;
2. we show that a navigation model trained purely on virtual observations performs poorly when tested on the corresponding real environment and how domain adaptation approaches can help to transfer the learned policy to the real world;

---

<sup>1</sup><https://matterport.com/>

3. we investigate the impact of sensor and actuator noise models on navigation performance by comparing navigation models trained with and without noise and tested in a noisy environment.

## 2 Related Work

Previous works have considered the problem of assessing the transfer of a navigation policy from a virtual environment to a real one. Specifically, The authors of [2] proposed to fine-tune navigation agents on a grid of real-world images acquired from specific locations. Differently from this approach, we rely on a much denser collection of real-world images to train and test our models, which is possible thanks to the proposed tool. The authors of [8, 9] recently proposed to perform paired experiments in simulation and reality. However, this requires access to a robotic platform in the training and testing stages. Our approach aims to move the analysis as much as possible to simulation. While research on this specific topic is still emerging, our work is related to other lines of research including embodied navigation simulators, visual navigation, and virtual-to-real domain adaptation.

### 2.1 Embodied Navigation Simulators

The development of advanced simulators [6, 7, 8, 12, 13, 4] together with the acquisition of large-scale virtual indoor environments [14, 15, 7, 16] have been instrumental to speed up the research on autonomous agents and give them the ability to navigate environments and interact with humans. These simulators allow to train navigation algorithms based on deep learning which would be unfeasible to train directly in real-world environments [3, 5, 17, 18].

Despite such advances, transferring a navigation policy trained in simulation directly to real contexts is not straightforward and still a problem under exploration [2, 8, 9]. Our work builds on previous work on simulators for visual navigation [4], aiming to augment them with realistic noise models and real-world observations.

### 2.2 Visual Navigation

In end-to-end visual navigation, the policy is learned by the agent given egocentric observations and the goal to reach. If the goal is given as the coordinates of the target or as an image, the task is referred to as *geometric navigation*. In this case, the agent is encouraged to learn the geometry of the environment in order to complete the task. Previous works have addressed geometric navigation in different settings, such as providing images of the goal as target for the agent [2], jointly learning the mapping of the environment and the policy to reach the destination in a supervised way [3], explicitly modelling environment mapping through SLAM and both long- and short-term goal estimation [5]. Other authors have shown that it is possible to exploit the efficiency of the Habitat simulator [4] to train a mapless PointGoal navigation policy with millions [19] or even billions [18] of frames, greatly outperforming classic approaches.

In this work, we consider the settings of geometric navigation in which the goal is provided in the form of coordinates relative to the agent’s position [19].

In our investigation, we focus on the agent’s ability to transfer the geometric representation learned in a virtual environment to its real counterpart.

### 2.3 Virtual-to-Real Domain Adaptation

Domain adaptation techniques can be employed to reduce the domain gap between virtual and real observations by learning domain-invariant image representations which allow to transfer action policies learned in simulation to the real world. Among the most notable approaches, domain randomization methods aim to improve generalization to real-world observations by increasing the size and variety of training data in tasks such as robot grasping [20] and visual navigation [21, 22]. Other approaches focus on learning domain invariant task-specific [23] or task-agnostic [24] higher order statistics of the scenes which are then used to train the agent’s policy. Recent works have used Generative Adversarial Networks (GANs) [25] to transform the appearance of images from a domain to another [11, 26, 27].

In this work, we assess the impact of simple adaptation techniques on the ability to transfer a policy learned with simulated observations to real-world observations. Specifically, we consider fine-tuning on real observations, which is made possible thanks to the proposed framework, and adapting virtual observations with CycleGAN [11], which is an unsupervised domain adaptation approach as it does not require labeled real images.

## 3 Embodied Visual Navigation In Real Environments Through Habitat

The proposed tool is based on the popular Habitat simulator [4], an open-source embodied AI 3D platform provided with flexible APIs that allow researchers to access and extend its functionalities. Habitat can render RGB, depth and semantic images at up to 10,000 fps and allows to import a variety of scenes, including custom environments acquired with a 3D scanner such as Matterport 3D. Our tool, as well as the dataset composed by 3D models and images used for evaluation purposes are publicly available at the following link: <https://iplab.dmi.unict.it/EmbodiedVN/>. As illustrated in Figure 2 and briefly discussed in the introduction, the proposed approach involves the acquisition of two 3D models of the environment, the alignment of the models, and the generation of navigation episodes using virtual and real-world observations, considering also sensor and actuator noise. We give details for each of these steps in the following sections.

### 3.1 Acquisition and Alignment of the 3D Models

As briefly mentioned in the introduction, the “virtual” 3D model is a geometrically accurate representation of the environment that can be imported in a simulator to perform efficient navigation episodes. The virtual model can be acquired using a scanner such as Matterport 3D, which is able to collect RGB and depth observations from different regions of the environment to reconstruct the 3D mesh of the space. The main limitation of the resulting 3D model is that its photo-realism is limited. The top part of Figure 1 compares a visual

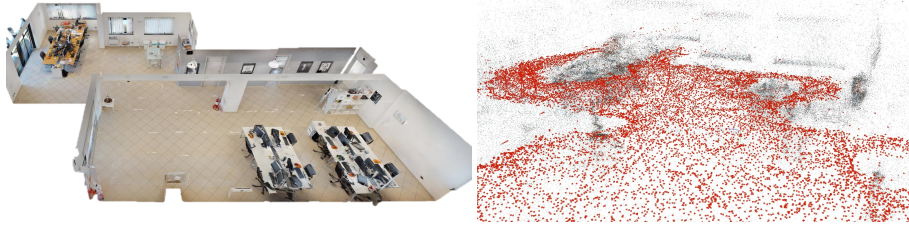


Figure 3: A view of the virtual (left) and real-world (right) 3D models. The virtual model is geometrically accurate and allows to sample images from any position but it lacks in photo-realism. The real-world model is a sparser collection of localized real-world images (each red marker represents the position of an image)

observation sampled from a virtual 3D model (on the left) to a corresponding real-world image (on the right).

The “real-world” 3D model consists in a sparse reconstruction of the environment which does not provide an accurate geometric representation of the scene but contains photo-realistic images. This model can be constructed starting from real images collected from the environment using a digital camera or a real robotic platform (e.g. with random navigation). The set of sparse RGB images can be used to obtain a sparse 3D point cloud of the environment in which each image is labeled with its camera pose, using a Structure from Motion (SfM) algorithm [10]<sup>2</sup>. Figure 3 compares the two models. We would like to note that the virtual 3D model can be loaded into the Habitat simulator using the official API, whereas the real-world 3D model cannot be easily interfaced with Habitat without the proposed tool.

Since the virtual and the real 3D models are constructed independently, they are characterized by different coordinate systems. We align them using an image-based registration procedure that allows to transform the coordinate system of the real-world 3D model to match the one of a reference set of images sampled from the same environment and annotated with their camera pose<sup>3</sup>. To obtain this reference set, we use the Habitat simulator to collect virtual images with random poses from the virtual 3D model. Despite the fact that the photo-realism of the virtual model is limited, we found the alignment procedure to succeed.

### 3.2 Generating Navigation Episodes with Real-World Observations

Once the virtual and the real 3D models are aligned, the proposed tool can produce navigation episodes with real-world observations. This is done by running the simulator on the virtual environment and systematically replacing virtual observations with the real-world which are closest in space to the current agent position. Specifically, at each navigation step, the current pose of the agent

<sup>2</sup>We use the COLMAP software: <https://colmap.github.io/faq.html>.

<sup>3</sup>We use the *model\_aligner* function of the COLMAP software <https://colmap.github.io/faq.html>



is used to perform a fast retrieval of the nearest image in the real-world 3D model. The retrieved real image can be fed by the visual navigation model in the simulator in order to choose the action to perform, together with the information about the goal to reach and other possible signals useful for the task. The retrieval process is repeated at each time-step until the navigation episode is completed.

To simplify retrieval, we transformed the 6DoF camera pose to 3DoF coordinates, where the first two degrees of freedom represent the X and the Z coordinates (parallel to the ground plane), whereas the third degree of freedom represents its rotation angle with respect to the Y axis perpendicular to the XZ plane. The transformation to 3DoF is performed because we assume that the agent’s altitude is constant during the navigation and that the agent’s camera does not perform rotations along its pitch and roll angles. To be able to retrieve angles using the cosine distance, we represent the camera rotation  $\theta$  as a unit vector  $(u, v) = (\cos \theta, \sin \theta)$ . We use the FAISS library [28] to perform an efficient two-steps retrieval procedure. In the first step, we filter images by angle, thresholding the cosine similarity score between the query and target angles. This allows us to select cameras which look in the same direction as the query pose. In the second step, we choose the image with the closest X-Z position among the filtered ones. In our experiments, we found that a cosine similarity threshold of 0.96 allows to obtain reasonable results.

### 3.3 Sensor and Actuator Noise Models

One of the major differences between simulation and reality is the presence of noise in perception and actions. As emerged in our experiments, navigation models trained in simulation assuming a perfect odometry perform poorly when deployed in noisy contexts and the problem can be alleviated by training them with simulated noise. We implemented two noise models: one for the sensor module (perception module), the other for the actuation module (action module). In our setup, the only sensor considered is a localization sensor which is used to determine the position of the agent at each time-step. Our model can work with generic noise models that can be tuned and adapted to any specific robotic platform. In this paper, we assume Gaussian noise for sensors and actuators.

## 4 Experimental Setup

In this section, we report details about the dataset acquisition, the simulation of the episodes, training and evaluation.

### 4.1 Environment and Simulation

For our experiments, we considered an office of about 150 square meters. We reconstructed the virtual and real-world 3D models of the same environment as described in Section 3. Specifically, the virtual 3D model has been acquired using Matterport 3D and imported in the Habitat simulator [4] to perform the virtual navigation episodes. The real 3D model has been created by collecting



Table 1: Noise levels considered in our experiments as standard deviations of our Gaussian noise models.

	Noise level		
	Small	Medium	Large
Localization noise	0.20m; 7°	0.40m; 15°	0.80m; 30°
Actuation noise	0.05m; 5°	0.10m; 10°	0.20m; 20°

24, 896 RGB images using the Sanbot Elf<sup>4</sup> robotic platform equipped with an Intel RealSense D435 camera<sup>5</sup>, following a simple random exploration procedure aimed at covering the whole space. The collected images have been used to reconstruct the 3D space with COLMAP [10]. With the reconstruction, we obtained as output the 6DoF camera pose of all images, relative to the reconstructed space. The set of real-world images has been split into training and test sets of equal size, following a sampling strategy to ensure a uniform distribution of camera poses in both sets and that neighboring frames do not fall into different sets to avoid over-fitting. We used the training set to train and fine-tune the navigation models on real-world images and the test set to evaluate their performances.

To align the 3D models, we used the aforementioned image-based alignment procedure offered by COLMAP [10], relying on a reference set of 6K images, labeled with their X, Y, Z camera coordinates. The images used for the alignment have been randomly sampled from the virtual environment using the Habitat simulator [4].

We trained our models following the setup proposed in [18], performing our navigation episodes in the Habitat simulator [4]. The navigation is formulated as a PointGoal navigation task [19]. At each time-step, the agent receives the current RGB observation and the updated goal coordinates relative to the current agent’s pose, then it chooses one of the four discrete actions to perform: *move straight by 0.25m*, *turn left by 10°*, *turn right by 10°*, *STOP*. The navigation episode ends when the *STOP* action is performed or when the maximum number of execution steps is reached. The maximum number of steps is set to 200 in our experiments.

We use the RGB DD-PPO Reinforcement Learning model [18] pre-trained on the Gibson [7] and Matterport 3D [15] datasets. The model is composed by a visual encoder and two recurrent layers. The visual encoder takes a  $256 \times 256$  RGB image as input and outputs an embedding, which is concatenated with the information about the goal coordinates and the previous action performed. The resulting vector is then fed to a recurrent network which outputs the action to be performed.

## 4.2 Domain Adaptation of the Navigation Policy to the Real Environment

We performed different experiments to evaluate the ability of the navigation model to adapt to real-world episodes. Specifically, we measured the perfor-

<sup>4</sup><http://en.sanbot.com/>

<sup>5</sup><https://www.intelrealsense.com/>

mance of the navigation model tested on real-world images when:

1. trained on virtual observations;
2. trained on real observations;
3. trained on virtual observations and fine-tuned on real observations;
4. trained on virtual observations adapted with CycleGAN [11];
5. trained on virtual observations and fine-tuned on virtual observations adapted with CycleGAN [11];
6. trained on virtual observations adapted with CycleGAN [11] and fine-tuned on real observations;
7. trained on virtual observations, fine-tuned on virtual observations adapted with CycleGAN [11] and then further fine-tuned on real observations.

CycleGAN has been considered in the experiments as a form of unsupervised domain adaptation to assess if it is beneficial to adapt the virtual observations to the real-world images. We trained CycleGAN [11] for 50 epochs using two sets of images: a set of 5k images randomly sampled from the virtual 3D model through Habitat [4] and another set of 5k images randomly sampled from the train set of the real-world images.

### 4.3 Navigation Episodes with Sensors and Actuation Noises

We considered three increasing levels of localization noise and three increasing levels of actuation noises, small, medium, and large, as reported in Table 1. The medium actuation noise values have been estimated from measurements collected performing actions such as “go straight by 0.25m” and “rotate by 10°” on the considered Sanbot Elf robotic platform. The medium localization noise values have been chosen considering the typical error of an image-based indoor localization system [9].

We evaluated the impact of localization and sensor noise by training visual navigation models in simulation with and without noise and testing them with progressive levels of noise. Specifically, we first trained the navigation models without noise and then fine-tuned them with small, medium and large sensor and actuator noise levels and with different combinations of them, after freezing the visual encoder weights.

### 4.4 Evaluation

We evaluated the performance of the visual navigation models in terms of SPL and success rate in reaching the goal in a limited number of steps, which are two standard evaluation measure for the visual navigation task. The SPL is a measure of the efficiency of the navigation episode with respect to the shortest geodesic path and is defined as:

$$\frac{1}{N} \sum_{i=1}^N S_i \frac{l_i}{\max(l_i, p_i)} \quad (1)$$

Table 2: Virtual to Real Policy Transfer

Training Stages	SPL	Success rate	Avg. dist. from goal (meters)	# training frames (Million)
Virtual	0.0160	0.022	7.9722	2.4
CycleGAN	0.2464	0.3310	4.6065	2.4
Virtual+CycleGAN	0.2648	0.3410	4.7535	1.2+1.2
Real	0.7112	0.8590	0.7709	2.4
Virtual+real	0.8001	0.9700	0.2493	1.2+1.2
CycleGAN+real	0.7665	0.8880	0.5219	1.2+1.2
Virtual+CycleGAN+real	0.7553	0.9360	0.3313	1.2+1.2+1.2

where  $N$  is the number of performed episodes,  $S_i$  is a boolean indicator of the success of the  $i$ -th episode,  $l_i$  is the shortest geodesic path length from the starting position to the goal position of the  $i$ -th episode and  $p_i$  is the agent’s path length in the  $i$ -th episode. If the navigation episode follows exactly the shortest path, it assumes the value of 1. On the contrary if the navigation policy fails, it assumes the value of 0. We randomly sampled 1000 navigation episodes defined by a starting and a goal positions, thus we stored them and used them to evaluate all the navigation models. To avoid to sample excessively simple navigation episodes, the navigation episodes have been filtered to ensure that the ratio between the geodesic distance and the euclidean distance from the starting position to the goal is greater than 1.1, as already suggested in [4]. We set the episodic maximum number of steps threshold to 200. We believe this is a reasonable value given the size and the complexity of the environment. An episode is considered successful if the agent calls the *STOP* action within  $0.20m$  from the goal and unsuccessful otherwise.

## 5 Results

### 5.1 Virtual to Real Policy Transfer

Table 2 reports the results of the compared methods when trained on different combinations of virtual and real observations in the absence of sensor and actuator noise, and tested on real-world observations. For each experiment, we report the set of training stages used to learn the policy. The notation  $A + B$  indicates that the model has been first trained on  $A$ , then fine-tuned on  $B$ . “Virtual” and “Real” denote that the model has been trained with virtual and real-world observations respectively, whereas “CycleGAN” indicates that the model has been trained on virtual observations translated to real using CycleGAN. The last column of the table reports the number of training frames seen in each stage. As can be seen, the approach trained only with virtual observations achieves very limited performance, e.g., “Virtual” obtains an SPL of 0.0160 and a success rate of 0.0220. Performance improves when CycleGAN is used to reduce the domain gap. Indeed, “CycleGAN” obtains a better success rate of 0.3310. Pre-training the model with  $1.2M$  frames on virtual observations, and then fine-tuning with observations translated with CycleGAN allows for minor improvements on all measures. For instance, the SPL of “Virtual + CycleGAN” is equal to 0.2648,

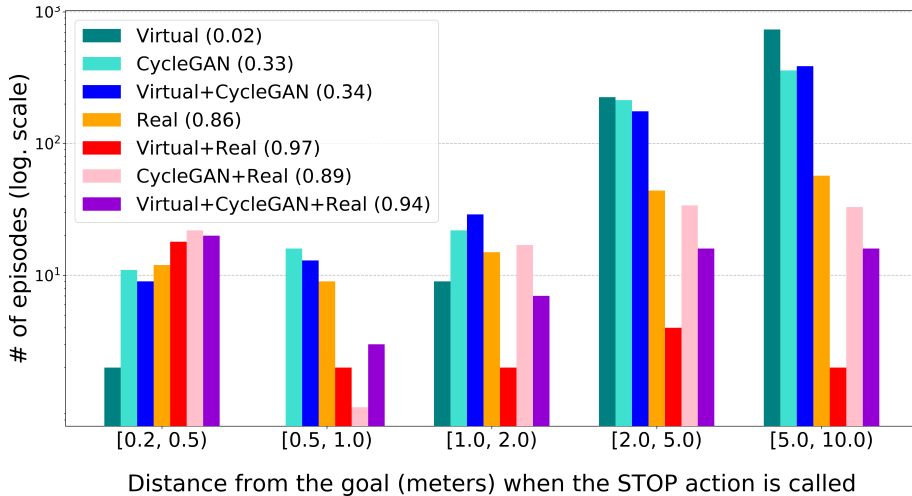


Figure 4: Distribution of unsuccessful navigation episodes with respect to the final distance from the goal for models evaluated on real-world observations. Success rates are reported in parenthesis in the legend.

which is better than “CycleGAN” by about two percentage points. Since the total number of frames seen during training by the two approaches is  $2.4M$  in both cases, we speculate that the minor improvement might be due to the virtual observations being less noisy than the ones translated with CycleGAN, which might improve generalization ability. It is worth noting that, while these improvements might seem scarce, the approaches using CycleGAN do not rely on any real-world labels such as the camera poses obtained using structure from motion. As such, they suggest a promising research direction for virtual to real policy transfer.

Interestingly, training the agent with real observations using the proposed tool allows to obtain major improvements in performance. For instance, “Real” achieves a success rate of 0.8590, which is a significant improvement over the success rate of 0.3410 obtained by “Virtual+CycleGAN”. Pre-training the model with virtual observations allows to obtain additional improvements bringing the SPL to 0.8001 (“Virtual + real”), the success rate to 0.97 and the average distance from the goal to 0.2493 meters. This result suggests that the knowledge learned in a purely simulated environment can be transferred to real-world observations using appropriate tools. Combining CycleGAN with training on real-world observations does not lead to improvements. We believe this is due to the limited ability of CycleGAN to translate virtual observations to real.

Figure 4 reports the distribution of the unsuccessful navigation episodes with respect to the final distance from the goal for all the methods discussed in this section. Specifically, each bin accumulates the number of episodes for which the distance of the agent from the goal when the STOP action has been called was in a given range. Since episodes are considered unsuccessful when this distance is above  $0.2m$ , the plot aims to assess qualitatively how far the agent is from the goal when it fails. As shown in the plot, in most of the unsuccessful episodes, the agent trained only on virtual observations (“Virtual”) tend to be located very far from the target with a peak in the range  $[5.0, 10.0)$ . Training with adapted

Table 3: Sensor and actuator noise with virtual observations

Sensors noise	Actuators noise	Trained with noise	SPL	Success rate	Avg. dist. from goal (meters)
No	No	No	0.9127	0.9910	0.1291
Small	No	No	0.8173	0.8910	0.1581
		Yes	0.8658	0.9380	0.1065
Medium	No	No	0.5075	0.5660	0.2404
		Yes	0.7114	0.7910	0.1554
Large	No	No	0.1552	0.1870	0.4909
		Yes	0.3643	0.4130	0.2577
No	Small	No	0.9092	0.9890	0.1171
		Yes	0.9073	0.9820	0.0956
No	Medium	No	0.8903	0.9700	0.1432
		Yes	0.8805	0.9740	0.1150
No	Large	No	0.8043	0.8860	0.2337
		Yes	0.8381	0.9340	0.2322
Small	Small	No	0.8020	0.8740	0.1876
		Yes	0.8328	0.8950	0.1607
Medium	Medium	No	0.4537	0.5100	0.2675
		Yes	0.4715	0.5290	0.2712
Large	Large	No	0.1288	0.1620	0.5442
		Yes	0.2450	0.2790	0.4040

visual observations (“CycleGAN” and “Virtual+CycleGAN”) allows to obtain a flatter distribution, with many failures moving to the range  $[0.2, 0.5)$ . When the model is trained with real observations, the failures are more evenly distributed across bins. Interestingly, “Virtual+Real” exhibits a distribution skewed towards the left side, which indicates that most of the unsuccessful episodes have been terminated in the proximity of the target. Combining CycleGAN with real data (“CycleGAN+Real” and “Virtual+CycleGAN+Real”) does not bring significant benefits as already observed in Table 2.

## 5.2 Influence of actuator and sensor noise

Table 3 reports the performance of navigation models when trained and tested in the presence of different combinations of sensor and actuator noise. This analysis is performed on virtual observations as these provide a much more accurate geometrical representation of the space. We consider two main approaches: 1) training the model without noise and testing it with noise, 2) training and testing the model with the same amount of noise. The visual navigation model trained and tested on virtual images in the absence of noise obtains good SPL and success rate values of 0.9127 and 0.9910 respectively (first row “No-No-No” in Table 3). This suggests that such methods can successfully learn near-optimal navigation policies when trained and tested in simulation. However, even adding small amounts of sensor or actuation noise during test degrades performance. For instance an SPL of 0.8658 is achieved when the model is tested with small sensor noise (“Small-No-No” in Table 3) and an SPL of 0.9073 when tested with

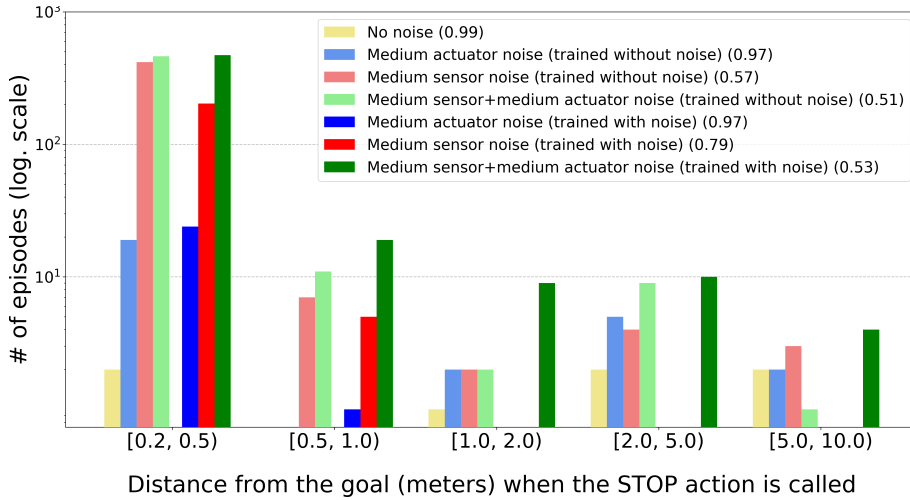


Figure 5: Distribution of unsuccessful navigation episodes with respect to the final distance from the goal for models trained and evaluated on virtual observations with and without sensor and actuator noise. Success rates are reported in parenthesis in the legend.

small actuator noise (“No-Small-No” in Table 3). The drop in performance is wider as larger amounts of noise and combinations of sensor and actuator noise are considered. For instance, the model obtains a small success rate of 0.1620 when tested with large sensor and actuator noise (“Large-Large-No” in Table 3). Interestingly, part of the gap in performance is generally recovered by training the methods with the same kind of noise. For instance, training the model with large amounts of sensor and actuator noise brings the success rate from 0.1620 to 0.2790 (“Large-Large-Yes” in Table 3). Similarly, training the model with small sensor noise bring the success rate from 0.8910 to 0.9380 (compare “Small-No-No” to “Small-No-Yes” in Table 3). Similar trends are observed with other combinations of sensor and actuator noise.

Figure 5 reports the distribution of the unsuccessful navigation episodes with respect to the final distance from the goal for some of the experiments considered in this section. As can be noted, while the distribution of unsuccessful episodes is overall uniform when the model is trained and tested without noise (“No noise” in Figure 5), methods trained without noise but tested with medium actuator or sensor noise tend to terminate unsuccessful episodes at large distances from the target (see “Medium sensor/actuator noise (trained without noise)” in Figure 5). The same approaches trained in the presence of noise (see “Medium sensor/actuator noise (trained with noise)” in Figure 5) are characterized by a distribution which is much more skewed towards small distances. This indicates that most of the episodes terminating at large distances from the target (e.g., at a distance larger than 1m) have been recovered, and most of the failures are due to the inability of the agent to reach the target in the final part of the navigation episode in which more precision is required. Different considerations apply to the case of the combination of medium sensor and actuator noise (compare “Medium sensor + medium actuator noise (trained without noise)” and

Table 4: Sensor and actuator noise with real-world observations

Sensors noise	Actuators noise	Trained with noise	SPL	Success rate	Avg. dist. from goal (meters)
No	No	No	0.8001	0.9700	0.2493
Small	Small	No	0.6553	0.7880	0.6564
		Yes	0.4708	0.6270	0.7269
Medium	Medium	No	0.4041	0.4041	0.7771
		Yes	0.3116	0.4150	0.8231
Large	Large	No	0.1467	0.1870	0.9808
		Yes	0.0827	0.1180	1.9109

“Medium sensor + medium actuator noise (trained with noise)” in Figure 5). Even if the success rate improves when the model is trained in the presence of noise, the two distributions do not show noticeable differences.

### 5.3 Influence of actuator and sensor noise with real-world observation

Table 4 finally reports the results of training and testing navigation policies on real-world observations in the presence of sensor and actuator noise. All models have been trained for 1.2M frames on virtual observations without noise, then fine-tuned for 1.2M frames on real-world observations with or without noise as specified in the table. Similarly to what observed in the case of virtual observations, models trained without noise perform worse when tested with noise. Indeed, while the baseline method has an SPL of 0.8001 (“No-No-No” in Table 4), the same model obtains an SPL of 0.6553 when trained with small sensor and actuator noise (“Small-Small-No” in Table 4), 0.4041 when trained with medium sensor and actuator noise (“Medium-Medium-No” in Table 4), and 0.1467 when trained with large sensor and actuator noise (“Medium-Medium-No” in Table 4). This confirms that noise plays an important role in the quality of the execution of navigation policies. Differently from what observed for virtual observations, training the same models with noise leads to even worse results (compare the performance of models trained with and without noise in Table 4). We speculate that this might be due to the more challenging nature of the navigation task with real-world observations and suggest that more sophisticated techniques to compensate for noise should be investigated.

## 6 Conclusion

We investigated the problem of transferring visual navigation policies trained in simulation to the real world. We proposed a tool based on Habitat [4] which can be used to produce virtual and real-world navigation episodes accounting for sensor and actuator noise models. Using the tool, we assessed the impact of the domain shift between virtual and real-world navigation induced by the difference in appearance and physical interactions. Our experiments suggest that adaptation methods are much needed to obtain visual navigation policies



able to generalize to the real world. We think that the proposed framework is a promising tool to assess and improve the generalization ability of navigation policies to the real world and hope that it will encourage research in this direction.

## Acknowledgment

This research is supported by OrangeDev s.r.l., by Piano della Ricerca 2016-2018 - CHANCE - Linea di Intervento 1 of DMI, University of Catania, and by MIUR AIM - Attrazione e Mobilità Internazionale Linea 1 - AIM1893589 - CUP E64118002540007.

## References

- [1] F. Bonin-Font, A. Ortiz, and G. Oliver, “Visual navigation for mobile robots: A survey,” *Journal of intelligent and robotic systems*, vol. 53, no. 3, p. 263, 2008.
- [2] 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 *international conference on robotics and automation (ICRA)*, 2017.
- [3] S. Gupta, J. Davidson, S. Levine, R. Sukthankar, and J. Malik, “Cognitive mapping and planning for visual navigation,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [4] Manolis Savva\*, Abhishek Kadian\*, Oleksandr 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 *International Conference on Computer Vision (ICCV)*, 2019.
- [5] D. S. Chaplot, D. Gandhi, S. Gupta, A. Gupta, and R. Salakhutdinov, “Learning to explore using active neural slam,” in *International Conference on Learning Representations (ICLR)*, 2020.
- [6] M. Savva, A. X. Chang, A. Dosovitskiy, T. Funkhouser, and V. Koltun, “MINOS: Multimodal indoor simulator for navigation in complex environments,” *arXiv:1712.03931*, 2017.
- [7] F. Xia, A. R. Zamir, Z.-Y. He, A. Sax, J. Malik, and S. Savarese, “Gibson env: real-world perception for embodied agents,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [8] M. Deitke, W. Han, A. Herrasti, A. Kembhavi, E. Kolve, R. Mottaghi, J. Salvador, D. Schwenk, E. VanderBilt, M. Wallingford, L. Weihs, M. Yatskar, and A. Farhadi, “RoboTHOR: An Open Simulation-to-Real Embodied AI Platform,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

- [9] A. Kadian, J. Truong, A. Gokaslan, A. Clegg, E. Wijmans, S. Lee, M. Savva, S. Chernova, and D. Batra, “Are we making real progress in simulated environments? measuring the sim2real gap in embodied visual navigation,” *International Conference on Intelligent Robots and Systems (IROS)*, 2020.
- [10] J. L. Schönberger and J.-M. Frahm, “Structure-from-motion revisited,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [11] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” in *International Conference on Computer Vision (ICCV)*, 2017.
- [12] Y. Wu, Y. Wu, G. Gkioxari, and Y. Tian, “Building generalizable agents with a realistic and rich 3d environment,” *arXiv preprint arXiv:1801.02209*, 2018.
- [13] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, D. Gordon, Y. Zhu, A. Gupta, and A. Farhadi, “Ai2-thor: An interactive 3d environment for visual ai,” *arXiv preprint arXiv:1712.05474*, 2017.
- [14] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma, A. Clarkson, M. Yan, B. Budge, Y. Yan, X. Pan, J. Yon, Y. Zou, K. Leon, N. Carter, J. Briales, T. Gillingham, E. Mueggler, L. Pesqueira, M. Savva, D. Batra, H. M. Strasdat, R. D. Nardi, M. Goesele, S. Lovegrove, and R. Newcombe, “The Replica dataset: A digital replica of indoor spaces,” *arXiv preprint arXiv:1906.05797*, 2019.
- [15] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, “Matterport3d: Learning from rgb-d data in indoor environments,” *International Conference on 3D Vision (3DV)*, 2017.
- [16] S. Song, F. Yu, A. Zeng, A. X. Chang, M. Savva, and T. Funkhouser, “Semantic scene completion from a single depth image,” *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [17] T. Chen, S. Gupta, and A. Gupta, “Learning exploration policies for navigation,” in *International Conference on Learning Representations (ICLR)*, 2019.
- [18] 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,” in *International Conference on Learning Representations (ICLR)*, 2020.
- [19] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva, *et al.*, “On evaluation of embodied navigation agents,” *arXiv preprint arXiv:1807.06757*, 2018.
- [20] 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 *International Conference on Intelligent Robots and Systems (IROS)*, 2017.

- [21] F. Sadeghi and S. Levine, “Cad2rl: Real single-image flight without a single real image,” *arXiv preprint arXiv:1611.04201*, 2016.
- [22] A. Loquercio, E. Kaufmann, R. Ranftl, A. Dosovitskiy, V. Koltun, and D. Scaramuzza, “Deep drone racing: From simulation to reality with domain randomization,” *IEEE Transactions on Robotics (T-RO)*, 2020.
- [23] Y. Zhang, P. David, and B. Gong, “Curriculum domain adaptation for semantic segmentation of urban scenes,” in *International Conference on Computer Vision (ICCV)*, 2017.
- [24] B. Sun, J. Feng, and K. Saenko, “Return of frustratingly easy domain adaptation,” in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [25] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *International Conference on Neural Information Processing Systems (NIPS)*, (Cambridge, MA, USA), 2014.
- [26] K. Rao, C. Harris, A. Irpan, S. Levine, J. Ibarz, and M. Khansari, “Rl-cyclegan: Reinforcement learning aware simulation-to-real,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [27] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, “Cycada: Cycle-consistent adversarial domain adaptation,” in *International conference on machine learning (ICML)*, 2018.
- [28] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with gpus,” *arXiv preprint arXiv:1702.08734*, 2017.