

RoboTHOR: An Open Simulation-to-Real Embodied AI Platform

Matt Deitke^{2*}, Winson Han^{1*}, Alvaro Herrasti^{1*}, Aniruddha Kembhavi^{1,2*}, Eric Kolve^{1*},
 Roozbeh Mottaghi^{1,2*}, Jordi Salvador^{1*}, Dustin Schwenk^{1*}, Eli VanderBilt^{1*}, Matthew Wallingford^{2*},
 Luca Weihs^{1*}, Mark Yatskar^{1*}, Ali Farhadi²

¹PRIOR @ Allen Institute for AI ²University of Washington

ai2thor.allenai.org/robothor

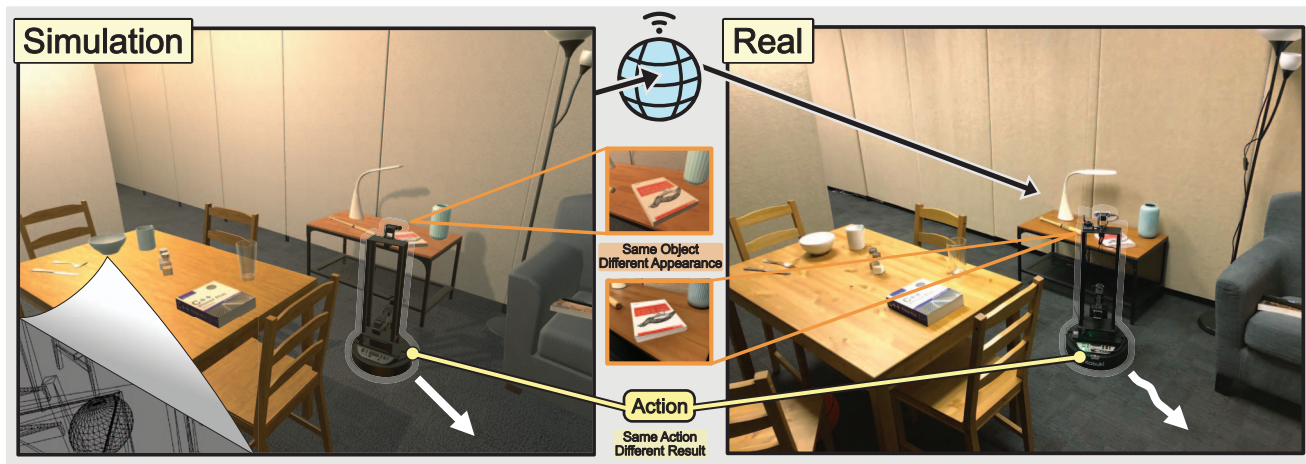


Figure 1: We present ROBOTHOR, a platform to develop and test embodied AI agents with corresponding environments in simulation and the physical world. The complexity of environments in ROBOTHOR along with disparities in appearance and control dynamics between simulation and reality pose new challenges and open many avenues for further research.

Abstract

Visual recognition ecosystems (e.g. ImageNet, Pascal, COCO) have undeniably played a prevailing role in the evolution of modern computer vision. We argue that interactive and embodied visual AI has reached a stage of development similar to visual recognition prior to the advent of these ecosystems. Recently, various synthetic environments have been introduced to facilitate research in embodied AI. Notwithstanding this progress, the crucial question of how well models trained in simulation generalize to reality has remained largely unanswered. The creation of a comparable ecosystem for simulation-to-real embodied AI presents many challenges: (1) the inherently interactive nature of the problem, (2) the need for tight alignments between real and simulated worlds, (3) the difficulty of replicating physical conditions for repeatable experiments, (4) and the associated cost. In this paper, we introduce ROBOTHOR to democratize research in interactive and embodied visual AI. ROBOTHOR offers a framework of simulated environments

paired with physical counterparts to systematically explore and overcome the challenges of simulation-to-real transfer, and a platform where researchers across the globe can remotely test their embodied models in the physical world. As a first benchmark, our experiments show there exists a significant gap between the performance of models trained in simulation when they are tested in both simulations and their carefully constructed physical analogs. We hope that ROBOTHOR will spur the next stage of evolution in embodied computer vision.

1. Introduction

For decades, the AI community has sought to create perceptive, communicative and collaborative agents that can augment human capabilities in real world tasks. While the advent of deep learning has led to remarkable breakthroughs in computer vision [33, 24, 48] and natural language processing [45, 16], creating active and intelligent embodied agents continues to be immensely challenging.

The widespread availability of large and open, computer

* Alphabetically listed equal contribution

vision and natural language datasets [50, 34, 47, 66], massive amounts of compute, and standardized benchmarks have been critical to this fast progress. In stark contrast, the considerable costs involved in acquiring physical robots and experimental environments, compounded by the lack of standardized benchmarks are proving to be principal hindrances towards progress in embodied AI. In addition, current state of the art supervised and reinforcement learning algorithms are data and time inefficient; impeding the training of embodied agents in the real world.

Recently, the vision community has leveraged progress in computer graphics and created a host of simulated perceptual environments such as AI2-THOR [32], Gibson [72], MINOS [53] and Habitat [54], with the promise of training models in simulation that can be deployed on robots in the physical world. These environments are free to use, continue to be improved and lower the barrier of entry to research in real world embodied AI; democratizing research in this direction. This has led to progress on a variety of tasks in simulation, including visual navigation [22, 69], instruction following [5, 68] and embodied question answering [21, 15]. But the elephant in the room remains: *How well do these models trained in simulation generalize to the real world?*

While progress has been ongoing, the large costs involved in undertaking this research has restricted pursuits in this direction to a small group of well resourced organizations. We believe that creating a free and accessible framework that pairs agents acting in simulated environments with robotic counterparts acting in the physical world will open up this important research topic to all—bringing faster progress and potential breakthroughs. As a step towards this goal, we present ROBOTHOR.

ROBOTHOR is a platform to develop artificial embodied agents in simulated environments and test them in both, simulation as well as the real world. A key promise of ROBOTHOR is to serve as an open and accessible benchmarking platform to stimulate reproducible research in embodied AI. With this in mind, it has been designed with the following properties:

- **Simulation and Real Counterparts** - ROBOTHOR consists of a training and validation corpus of 75 scenes in simulation at present. A total of 14 test-dev and test-standard scenes are present in simulation and their counterparts constructed in the physical world. Scenes are designed to have diverse wall and furniture layouts, and all are densely populated by a variety of object categories. Figure 1 shows a view of the dining room in one of the test-dev scenes, in simulation as well as real.
- **Modular** - Scenes in ROBOTHOR are built in a modular fashion, drawing from an asset library containing

wall structures, flooring components, ceilings, lighting elements, furniture pieces and objects; altogether totaling 731 unique assets distributed across scenes. This enables scene augmentation and easy expansion of ROBOTHOR to fit the needs of researchers.

- **Re-configurable** - The physical environments are also built using modular and movable components, allowing us to host scenes with vastly different layouts and furniture arrangements within a single physical space. This allows us to scale our test corpora while limiting their cost and physical footprint. Reconfiguring the space to a new scene can be accomplished in roughly 30 minutes.
- **Accessible to all** - The simulation environment, assets and algorithms we develop will be open source. More critically, researchers from all over the world will be able to remotely deploy their models on our hardware at no cost to them. We will be setting up a systematic means of reserving time in our environment.
- **Replicable** - The physical space has been designed to be easily replicable by other researchers should they wish to construct their own physical environment. This is achieved through open sourced plans, readily available building components, IKEA furniture, and a low cost amounting to roughly \$10,000 in materials and assets to create the physical space. In addition, we use LoCoBot, an inexpensive and easily obtainable robot.
- **Benchmarked** - In addition to open sourcing baseline models, we will host challenges involving several embodied AI tasks with a focus on the ability to transfer these models successfully onto robots running in a physical environment.

ROBOTHOR has been designed to support a variety of embodied AI tasks. In this work we benchmark models for semantic navigation, the task of navigating to an instance of the specified category in the environment. The complexity and density of scenes in ROBOTHOR renders this task quite challenging, with humans requiring 49.5 steps (median statistic) to find the target object. We train a set of competitive models using a pure reinforcement learning (RL) approach with asynchronous actor critic (A3C) on the simulated training environments, measure their performance on the simulated as well as real validation environments and arrive at the following revealing findings. (1) Similar to findings in past works such as [69], semantic navigation models struggle with generalizing to unseen environments in simulation. We show that their performance takes an even larger hit when deployed onto a physical robot in the real world. (2) We analyze simulated and real world egocentric views and find a disparity in feature space in spite of the images from the two modalities looking fairly similar to the naked

eye; a key factor affecting the transfer of policies to the real world. (3) As expected and noted in previous works, control dynamics in the real world vary significantly owing to motor noise, slippage, and collisions. (4) Off the shelf image correction mechanisms such as image-to-image translation do not improve performance.

These findings reveal that training embodied AI models that generalize to unseen simulated environments and further yet to the real world remains a daunting challenge; but also open up exciting research frontiers. We hope that ROBOTHOR will allow more research teams from across the globe to participate in this research which will result in new model architectures and learning paradigms that can only benefit the field.

2. Related Work

Embodied AI Environments. In recent years, several synthetic frameworks have been proposed to investigate tasks including visual navigation, task completion and question answering in indoor scenes [32, 53, 70, 8, 74, 72, 46, 54]. These free virtual environments provide excellent testbeds for embodied AI research by abstracting away the noise in low-level control, manipulation and appearance and allowing models to focus on the high-level end goal. ROBOTHOR provides a framework for studying these problems as well as for addressing the next frontier: transferring models from simulation to the real world.

Robotics research platforms [1, 2] have traditionally been expensive to acquire. More recent efforts have led to low cost robot solutions [58, 3, 43] opening up the space to more research entities. There has been a long history of using simulators in conjunction with physical robots. These largely address tasks such as object manipulation using robotic arms [11] and autonomous vehicles [58, 55].

Visual Navigation. In this paper, we explore models for the task of visual navigation, a popular topic in the robotics and computer vision communities. The navigation problem can be divided into two broad categories, spatial navigation and semantic navigation. Spatial navigation approaches [62, 17, 64, 23, 25, 56, 79, 49, 30, 14, 10] typically address navigating towards a pre-specified coordinate or a frame of a scene and they focus on understanding the geometry of the scene and learning better exploration strategies. For example, [79] address navigation towards a given input image, [10] address navigation towards a point in a scene and [30] learn a collision-free navigation policy. Semantic navigation approaches [22, 75, 69, 71, 39, 52, 41] attempt to learn the semantics of the target in conjunction with navigation. For example, [75] use prior knowledge of object relations to learn a policy that better generalize to unseen scenes or objects. [69] use meta-learning to learn a self-supervised navigation policy toward a specified object category. [71] use prior knowledge of scene layouts to navigate to a spe-

cific type of room. We benchmark models on the task of semantic navigation.

Navigation using language instructions has been explored by [5, 68, 18, 31, 67, 36, 37]. This line of work has primarily been tested in simulation; transferability to the real world remains an open question and addressing this via ROBOTHOR is a promising future endeavour. Navigation has also been explored in other contexts such as autonomous driving (e.g., [12, 73]) or city navigation (e.g., [38, 13]). In this work, we focus on indoor navigation.

Sim2Real Transfer. Domain adaptation in general as well as Sim2Real in particular, have a long history in computer vision. There are different techniques to adapt models from a source domain to a target domain. The main approaches are based on randomization of the source domain to better generalize to the target domain [63, 29, 51, 44, 61], learning the mapping between some abstraction or higher order statistics of the source and target domains [27, 59, 19, 35, 76, 42, 78], interpolating between the source and the target domain on a learned manifold [20, 9], or generating the target domain using generative adversarial training [7, 60, 57, 6, 26, 28]. ROBOTHOR enables source randomization via scene diversity and asset diversity. We also experiment with using an off the shelf target domain mapping method, the GAN-based model of [77].

3. RoboTHOR

State of the art learning algorithms for embodied AI use reinforcement learning based approaches to train models, which typically require millions of iterations to converge to a reasonable policy. Training policies in the real world with real robots would take years to complete, due to the mechanical constraints of robots. Synthetic environments, on the other hand, provide a suitable platform for such training strategies, but how well models trained in simulation transfer to the real world, remains an open question. ROBOTHOR is a platform, built upon the AI2-THOR framework [32] to build and test embodied agents with an emphasis on studying this problem of domain transfer from simulation to the real world.

Scenes. ROBOTHOR consists of a set of 89 apartments, 75 in train/val (we use 60 for training and 15 for validation), 4 in test-dev (which are used for validation in the real world) and 10 in test-standard (blind physical test set) drawn from a set of 15, 2 and 5 wall layouts respectively. Apartments that share the same wall layout have completely different room assignments and furniture placements (for example, a bedroom in one apartment might be an office in another). Apartment layouts were designed to encompass a wide variety of realistic living spaces. Figure 3 shows a heatmap of wall placements across the train/val subset. This set of apartments is only instantiated in simulation, while the test-dev and test-standard apartments are also substantiated in

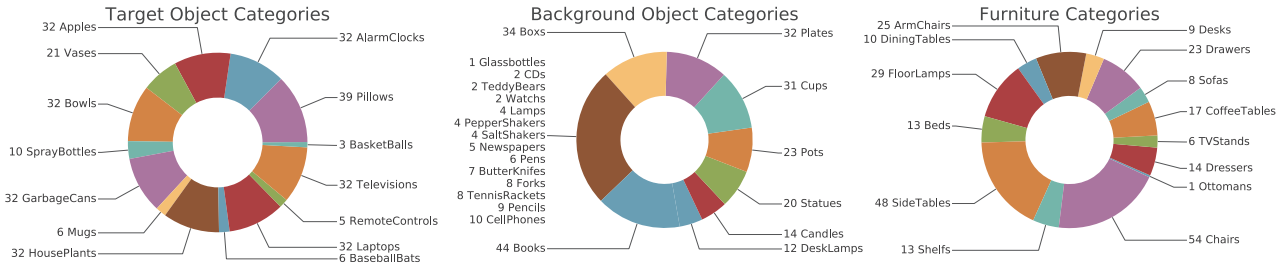


Figure 2: Distribution of object categories in ROBOTHR

the physical world. The layouts, furniture, objects, lighting, etc. of the simulation environments have been designed carefully so as to closely resemble the corresponding scenes in their physical counterparts, while avoiding any overlap between the wall layouts and object instances among train/val, test-dev and test-standard. This resemblance will enable researchers to study the discrepancies between the two modalities and systematically identify the challenges of the domain transfer.

Assets. A guiding design principle of ROBOTHR is *modularity*, which allows us to easily augment and scale scenes. A large asset library was created by digital artists from which scenes were created by selectively drawing from these assets. This is in contrast to environments that are based on 3D scans of rooms which are challenging to alter and interact with. The framework includes 11 types of furniture (e.g. TV stands and dining tables) and 32 types of small objects (e.g. mugs and laptops) across all scenes. The majority of real furniture and objects were gathered from IKEA. Among the small objects categories, 14 are designated as targets and guaranteed to be found in all scenes for use in semantic navigation tasks. In total there are 731 unique object instances in the asset library with no overlap among train/val, test-dev and test-standard scenes. Figure 2 shows the distribution of object categories amongst the asset library. We distribute object categories as uniformly as possible to avoid bias toward specific locations. Figure 3 shows the spatial distribution of target objects, background objects and furniture in the scenes. Figure 4 shows the distribution of the number of visible objects in a single frame. A large number of frames consist of the agent looking at the wall as is common in apartments, but outside these views many objects are visible to the agent at any given point as it navigates the environment.

Physical space. The physical space for ROBOTHR is $8.8m \times 3.9m$. The space is partitioned into rooms and corridors using ProPanel walls, which are designed to be lightweight and easy to set up and tear down, allowing us to easily configure a new apartment layout in a few minutes.

Agent. The physical robot used is a LoCoBot¹, which

¹<http://www.locobot.org/>

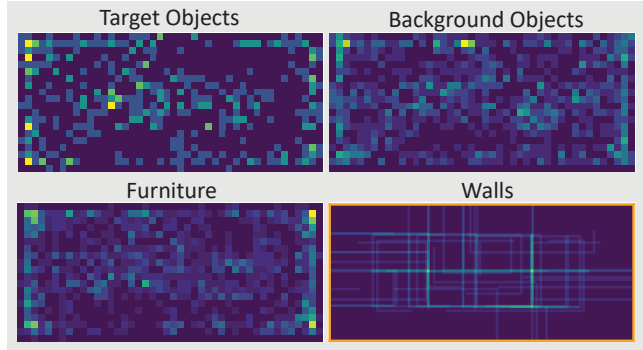


Figure 3: **Spatial distribution of objects and walls.** Heatmaps illustrate the diverse spatial distribution of target objects, background objects, furniture, and walls.

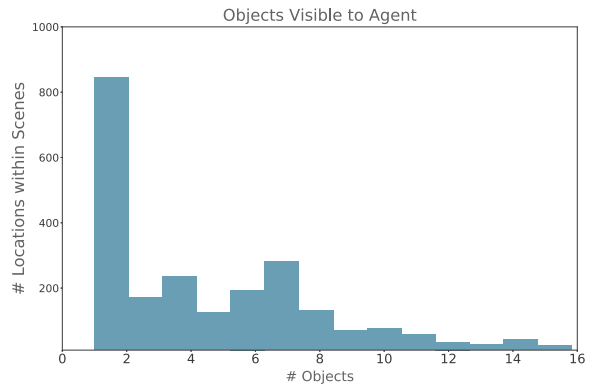


Figure 4: **Object visibility statistics.** The distribution of objects visible to an agent at a single time instant.

is equipped with an Intel RealSense RGB-D camera. We replicate the robot in simulation with the same physical and camera properties. To mimic the noisy dynamics of the robot movements, we add noise to the controller in simulation. The noise parameters are estimated by manually measuring the error in orientation and displacement over multiple runs.

API. To enable a seamless switch between the synthetic and the real environments, we provide an API that is agnostic to the underlying platform. Hence, agents trained in sim-

ulation can be easily deployed onto the LoCoBot for testing. The API was built upon the PyRobot [43] framework to manage control of the LoCoBot base as well as camera.

Connectivity. A main goal of this framework is to provide access to researchers across the globe to deploy their models onto this physical environment. With this in mind, we developed the infrastructure for connecting to the physical robot or the simulated agent via HTTP. A scheduler prevents accessing the same physical hardware by multiple parties.

Localization. We installed Super-NIA-3D localization modules across the physical environment to estimate the location of the robot and return that to the users. For our experiments, we do not use the location information for training as this type of training signal is not usually available in real world scenes. We use this location information only for evaluation and visualization.

4. Visual Semantic Navigation

In this paper, we benchmark models for the task of visual semantic navigation, i.e. navigating towards an instance of a pre-specified category. ROBOTHOR enables various embodied tasks such as question answering, task completion and instruction following. Navigation is a key component of all these tasks and is a necessary and important first step towards studying transfer in the context of embodied tasks.

Visual semantic navigation evaluates the agent’s capabilities not only in avoiding obstacles and making the right moves towards the target, but also understanding the semantics of the scene and targets. The agent should learn how different instances of an object category look like and should be able to reason about occlusion, scale changes and other variations in object appearance.

More specifically, our goal is to navigate towards an instance of an object category specified by a noun (e.g., Apple) given ego-centric sensory inputs. The sensory input can be an RGB image, a depth image, or combination of both. At each time step the agent must issue one of the following actions: *Move Ahead*, *Rotate Right*, *Rotate Left*, *Look Up*, *Look Down*, *Done*. The action *Done* signifies that the agent reports that it has reached its goal and leads to an end of episode. We consider an episode successful if (a) the object is in view (b) the agent is within a threshold of distance to the target and (c) the agent reports that it observes the object. The starting location of the agent is a random location in the scene.

The motion of the agent in the simulated world is stochastic in nature, mirroring its behavior in the real world. This renders the task more challenging. Previous works such as [69] consider agent motion along the axes on a grid. But given the end goal of navigating in the real world with motor noise and wheel slippage, deterministic movements in training lead to sub optimal performance during testing.

The semantic navigation task is very challenging owing

to the size and complexity of the scenes. Figure 5 shows the lengths of shortest paths to the target objects, in terms of the *Move Ahead* and *Rotate* actions. But shortest path statistics are a bit misleading for the task of semantic navigation because they assume that the agent already knows the location it must travel towards. In fact, the agent must explore until it observes the target, and then move swiftly towards it. We conducted a study where humans were posed with the problem of navigating in scenes in ROBOTHOR (simulation) to find targets. The median number of steps was 49.5 (compared to 22.0 for shortest paths), illustrating the exploration nature of the task. Figure 6 shows an example trajectory from a human compared to the corresponding shortest path.

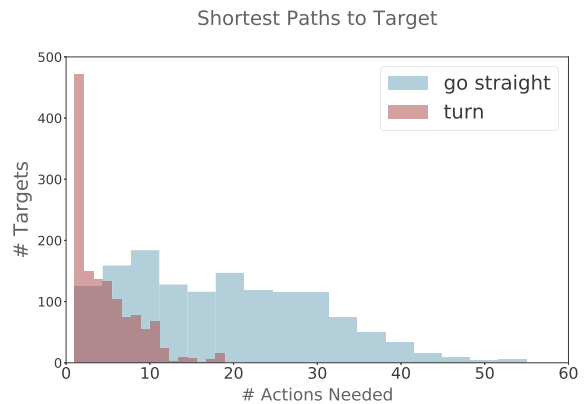


Figure 5: **Histogram of actions along the shortest path.** The number of actions invoked along the shortest paths to targets in the training scenes. Note that the shortest path is very difficult to obtain in practice, since it assumes *a priori* knowledge of the scene.



Figure 6: **Example human trajectory.** The shortest path to a target vs the path taken by a human from the same starting location are visualized. The human wanders around looking for the TV and on seeing it, walks straight towards it.

4.1. Baseline models

We measure the performance of the following baseline models:

Random - This model chooses an action randomly amongst the set of possible actions. This includes invoking the *Done* action. The purpose of this baseline is to ascertain if the scenes and starting locations are not overly simplistic.

Instant Done - This model invokes the *Done* action at the very first time step. The purpose of this baseline is to measure the percentage of trivial starting locations.

Blind - This model receives no sensory input. It only consists of an LSTM with simply a target embedding input. Its purpose is to establish a baseline that can only leverage starting and target location bias in the dataset.

Image-A3C - The agent perceives the scene at time t in the form of an image o_t . The image is fed into a pre-trained and frozen ResNet-18 to obtain a $7 \times 7 \times 512$ tensor, followed by two 1×1 convolution layers to reduce the channel depth to 32 and finally concatenated with a 64 dimensional embedding of the target word and provided to an LSTM which generates the policy. The agent is trained using the asynchronous advantage actor-critic (A3C) [40] formulation, to act to maximize its expected discounted cumulative reward. Two kinds of rewards are provided: a positive success reward and a negative step reward to encourage efficient paths.

Image+Detection-A3C - The image is fed into FasterRCNN [48] trained on the MSCOCO dataset. Each resulting detection has a category, probability and box locations and dimensions, which are converted to an embedding using an MLP. The resulting set of embeddings are converted into a tensor $7 \times 7 \times 120$ by aligning detection boxes to spatial grid cells and only considering the top 3 boxes per location. This tensor is concatenated to the tensor obtained for the image modality and then processed as above.

4.2. Metrics

We report Success Rate and Success weighted by Path Length (SPL) [4], common metrics for evaluating navigation models. In addition, we also report path lengths in terms of the number of actions taken as well as distance travelled. This quantifies the exploration carried out by the agent.

5. Experiments

Training. We train models to navigate towards a single target, a Television with an action space of 4 including *Move Ahead*, *Rotate Right*, *Rotate Left* and *Done*. The rotation actions rotate the agent by 45 degrees. The agent must issue the *Done* action when it observes the object. If an agent is successful, it receives a reward of +5. The step penalty at each time step is -0.01. To mimic the noisy dynamics of the real robot, we add noise to the movement of the virtual agent. For translation, we add a gaussian noise with mean 0.001m and standard deviation 0.005m, and for rotation we use a mean of 0° and standard deviation of 0.5° .

Models were trained on 8 TITAN X GPUs for 100,000 episodes using A3C with 32 threads. For each episode, we sample a random training scene and random starting location in the scene for the agent. We use the Adam optimizer with a learning rate of 0.0001.

We train the models on a subset of 50 train scenes and report numbers on 2 of the test-dev scenes. We report metrics on starting locations categorized into *easy*, *medium* and *hard*. If the length of the shortest path from the starting point to the target is among the lowest 20% path lengths in that scene, it is considered easy. If it is between 20% and 60%, it is considered medium, and longer paths are considered as hard.

Sim-to-Sim Table 1 shows the results of our benchmarked models when trained on simulation and evaluated on the test-dev scenes in simulation. The trivial baselines Random, Instant Done and Blind perform very poorly, indicating that the dataset lacks a bias that is trivial to exploit using these models. The Image only model performs reasonably well, succeeding at over 50% of easy trajectories but doing very poorly at hard ones. Adding object detection does not help performance. Our analysis in Section 6 shows that object detectors trained in the real world show a drop in performance in simulation, which might be contributing to their ineffectiveness.

Sim-to-Real Due to the slow motion of the robot, episodes on the robot last as long as 10 minutes. This limits the amount of testing that can be done in the real world. Table 2 shows the result of the best performing model on Sim-to-Sim, evaluated on a real scene on a subset of starting locations as those reported in Table 1. This shows that there is a sizeable drop in performance for the real robot, especially in SPL. The robot however, does learn to navigate around and explore its environment fairly safely and over 80% of the trajectories have no collisions with obstacles. This leads to high values for episode lengths in all 3 cases - easy, medium and hard.

Overfit Sim-to-Real The semantic navigation sim-to-real task in ROBOTHOR tests two kinds of generalization: Moving to new scenes and moving from simulation to real. To factor out the former and focus on the latter, we trained a policy on a test-dev scene, (which expectedly led to overfitting on sim) and then deployed this on the robot. Table 3 shows these results and demonstrates the upper bound of current models in the real world if they had memorized the test-dev scene perfectly in simulation. The robot does very well on the easy targets, but is affected a lot on hard targets. For all three modes, the SPL is affected tremendously. Appearance and control variations often lead the robot to spaces away from the target, which does not happen in simulation due to overfitting.

	Easy				Medium				Hard			
	Success	SPL	Episode length	Path length	Success	SPL	Episode length	Path length	Success	SPL	Episode length	Path length
Random	7.58	5.32	4.36	0.34	0.00	0.00	4.27	0.30	0.00	0.00	3.06	0.19
Instant Done	4.55	3.79	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00
Blind	4.55	3.79	1.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	1.00	0.00
Image	55.30	38.12	45.87	9.26	28.79	19.12	78.49	14.82	1.47	0.97	81.09	14.22
Image+Detection	36.36	19.89	63.41	11.39	11.36	5.25	90.37	16.65	0.74	0.61	83.01	14.00

Table 1: **Benchmark results for Sim-to-Sim**

	Easy				Medium				Hard			
	Success	SPL	Episode length	Path length	Success	SPL	Episode length	Path length	Success	SPL	Episode length	Path length
Image	33.33	3.53	53.16	7.18	16.66	3.70	43.83	5.33	0.00	0.00	67.83	7.00

Table 2: **Benchmark results for Sim-to-Real**

6. Analysis

We now dig deeper into the appearance disparities between real and simulation images via t-sne embeddings, object detection results and the output policies for both modalities. We also provide a study showing the effect of changing camera parameters between real and simulation for the transfer problem. Finally, we evaluate using an image translation method for the purposes of domain adaptation.

Appearance disparities. To the naked eye, images between the real and simulation worlds in ROBOTHOR look quite similar. However when we look at embeddings provided by networks, the disparity becomes more visible. We considered 846 images, each from simulation and real collected from the same locations in the scene, passed these images through ResNet-18 to obtain a 512 dim feature vector and then used t-SNE [65] to reduce the dimensionality to 3. Figure 7 shows these embeddings. One can clearly see the separation of points into real and simulation clusters. This indicates that embeddings of images in the two modalities are different, which goes towards explaining the drop in performance between simulation and real test-dev scenes (see Table 1 and Table 2). Figure 7 also shows the nearest neighbor (cosine similarity) simulation image to one of the real images. We found that although nearest neighbors are not far away spatially, they are still slightly different, sometimes having a different view, or a different distance to an obstacle. This might explain why the robot takes different actions in the real world. This analysis indicates that methods of representation learning might need to be revisited, especially when these representations must work across real and simulation modalities.

Object detection transfer. Since we leverage an off the shelf object detector (Faster-RCNN) trained on natural images (MS-COCO), it is imperative to compare the accuracy of this model on images collected in the real and simulated

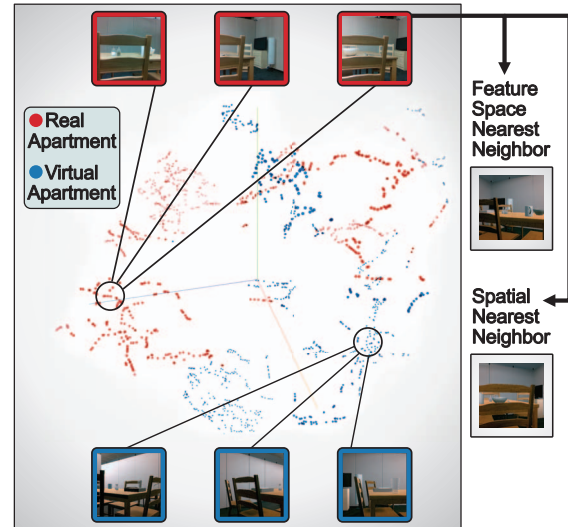


Figure 7: **Comparison of embeddings for real and synthetic images.** The scatter plot shows a t-SNE visualization of ResNet-18 (pre-trained on ImageNet) features for images from real and simulated apartments. Also shown are the nearest neighbor in feature space and spatial nearest neighbor, which differ slightly in the viewpoint of the agent.

apartments. We collected 761 corresponding images from both environments, ran Faster-RCNN on them and obtained ground truth annotations for 10 object classes (all of which intersected with MS-COCO classes) on Amazon Mechanical Turk. At an Intersection-over-Union (IOU) of 0.5, we obtained a mAP of 0.338 for real images and 0.255 for simulated ones; demonstrating that there is a performance hit going across modalities. Furthermore, detection probabilities tend to differ between the two modalities as demonstrated in Figure 8, rendering transfer more challenging for

	Easy				Medium				Hard			
	Success	SPL	Episode length	Path length	Success	SPL	Episode length	Path length	Success	SPL	Episode length	Path length
Image Sim-2-Sim	100	82.17	8.09	1.05	100	86.17	27.52	4.77	94.12	83.18	42.26	7.90
Image Sim-2-Real	100	12.28	15.00	2.33	83.33	18.68	43.33	5.5	50	28.53	30.16	7.54

Table 3: **Benchmark results for Sim-to-Real trained on a single test-dev scene.**

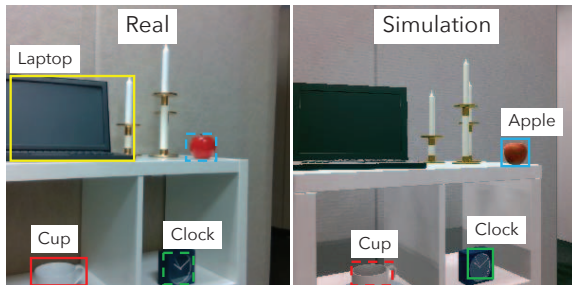


Figure 8: **Object detection.** Results of object detection in a real and simulated image. Solid lines denote high confidence detections whereas dashed lines denote low confidence detections.

models that exploit these probabilities. Note that this transfer is in the opposite direction (Real to Sim) compared to our current transfer setup, but is revealing nonetheless.

Modifying camera parameters. Since the real and simulation apartments were both designed by us, we were able to model the simulation camera close to the one present on LoCoBot. However, to test the sensitivity of the learnt policy to the camera parameters, we performed an experiment where we trained with a Field Of View (FOV) of 90° and tested with an FOV of 42.5° . The resulting real world experiments for the Image only model show a huge drop of performance to 16% for easy, and 0% for medium and hard. Interestingly, the robot tends to move very little and instead rotate endlessly. We hypothesize that models trained on simulation likely overfit to the image distribution observed at training, and the images captured at a different FOV vary so significantly in feature space, that they invoke unexpected behaviors. Since popular image representation models are trained on images from the internet, they are biased towards the distribution of cameras that people use to take pictures. Cameras on robots are usually quite different. This suggests that we should fine tune representations to match robot cameras and also consider different camera parameters as a camera augmentation step during training in simulation.

Domain adaptation via image translation. Since appearance statistics vary between real and simulation, we experimented with applying an image-to-image translation technique, CycleGAN [77] to translate real world images towards simulation images. This would enable us to train in simulation and apply the policy on the real robot while processing the translated images. We needed to use a trans-

lation model that could use unpaired image data, since we are unable to obtain paired images with 0 error in the placement of the agent. Paired images were obtained for 3 test-dev scenes to train CycleGAN. The policy (trained in simulation) was run on the robot on the remaining test-dev scene. Interestingly, the CycleGAN model does learn to flatten textures and adjust lighting and shadows as seen in Figure 9. However, the resultant robot performance is very poor and obtains 0% accuracy. While image translation looks pleasing to the eye, it does introduce spurious errors which hugely affect the image embeddings and thus the resultant policy.

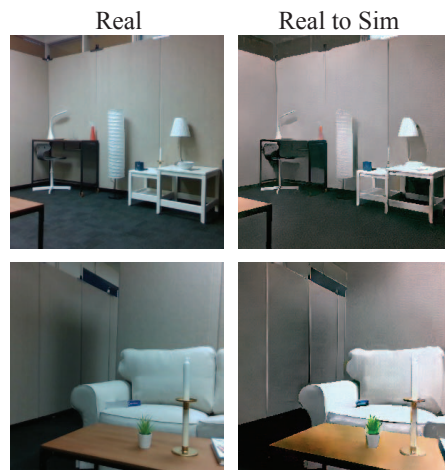


Figure 9: **Examples of real to simulation transfer.** We use CycleGAN [77] to translate real images towards simulated ones. The model learns to flatten out the texture and adjust the shadows to look like a simulated image.

7. Conclusion

In this paper, we presented ROBOTHOR, an open, modular, re-configurable and replicable embodied AI platform with counterparts in simulation and the real world, where researchers across the globe can remotely deploy their models onto physical robots and test their algorithms in the physical world. Our preliminary findings show the performance of models drops significantly when transferring from simulation to real. We hope that ROBOTHOR will enable more research towards this important problem.

Acknowledgements. This work is in part supported by NSF IIS 1652052, IIS 17303166, DARPA N66001-19-2-4031, 67102239 and gifts from Allen Institute for AI.

References

- [1] Baxter. [https://en.wikipedia.org/wiki/Baxter_\(robot\)](https://en.wikipedia.org/wiki/Baxter_(robot)). 3
- [2] Sawyer. <https://www.rethinkrobotics.com/sawyer>. 3
- [3] The MIT RACECAR. <https://mit-racecar.github.io>, 2016. 3
- [4] Peter Anderson, Angel X. Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, and Amir Roshan Zamir. On evaluation of embodied navigation agents. *arXiv*, 2018. 6
- [5] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. Vision-and-Language Navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*, 2018. 2, 3
- [6] Konstantinos Bousmalis, Alex Irpan, Paul Wohlhart, Yunfei Bai, Matthew Kelcey, Mrinal Kalakrishnan, Laura Downs, Julian Ibarz, Peter Pastor, Kurt Konolige, Sergey Levine, and Vincent Vanhoucke. Using simulation and domain adaptation to improve efficiency of deep robotic grasping. In *ICRA*, 2018. 3
- [7] Konstantinos Bousmalis, Nathan Silberman, David Dohan, Dumitru Erhan, and Dilip Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *CVPR*, 2017. 3
- [8] Simon Brodeur, Ethan Perez, Ankesh Anand, Florian Golemo, Luca Celotti, Florian Strub, Jean Rouat, Hugo Larochelle, and Aaron C. Courville. Home: a household multimodal environment. *arXiv*, 2017. 3
- [9] Rui Caseiro, Joao F. Henriques, Pedro Martins, and Jorge Batista. Beyond the shortest path : Unsupervised domain adaptation by sampling subspaces along the spline flow. In *CVPR*, 2015. 3
- [10] Devendra Singh Chaplot, Saurabh Gupta, Dhiraj Gandhi, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural mapping. In *CVPR Workshop on Habitat Embodied Agents*, 2019. 3
- [11] Yevgen Chebotar, Ankur Handa, Viktor Makoviychuk, Miles Macklin, Jan Issac, Nathan D. Ratliff, and Dieter Fox. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. In *ICRA*, 2019. 3
- [12] Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. In *ICCV*, 2015. 3
- [13] Howard Chen, Alane Suhr, Dipendra Misra, Noah Snaveley, and Yoav Artzi. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In *CVPR*, 2019. 3
- [14] Tao Chen, Saurabh Gupta, and Abhinav Gupta. Learning exploration policies for navigation. In *ICLR*, 2019. 3
- [15] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *CVPR*, 2018. 2
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. 1
- [17] Dieter Fox, Wolfram Burgard, and Sebastian Thrun. Markov localization for reliable robot navigation and people detection. In *Sensor Based Intelligent Robots*, 1999. 3
- [18] Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. Speaker-follower models for vision-and-language navigation. In *NeurIPS*, 2018. 3
- [19] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 2016. 3
- [20] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *CVPR*, 2012. 3
- [21] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. IQA: Visual question answering in interactive environments. In *CVPR*, 2018. 2
- [22] Saurabh Gupta, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. Cognitive mapping and planning for visual navigation. In *CVPR*, 2017. 2, 3
- [23] Raia Hadsell, Pierre Sermanet, Jan Ben, Ayse Erkan, Marco Scoffier, Koray Kavukcuoglu, Urs Muller, and Yann LeCun. Learning long-range vision for autonomous off-road driving. *J. of Field Robotics*, 2009. 3
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1
- [25] Peter Henry, Christian Vollmer, Brian Ferris, and Dieter Fox. Learning to navigate through crowded environments. In *ICRA*, 2010. 3
- [26] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. CyCADA: Cycle-consistent adversarial domain adaptation. In *ICML*, 2018. 3
- [27] Judy Hoffman, Dequan Wang, Fisher Yu, and Trevor Darrell. Fcns in the wild: Pixel-level adversarial and constraint-based adaptation. *arXiv*, 2016. 3
- [28] Haoshuo Huang, Qixing Huang, and Philipp Krahenbuhl. Domain transfer through deep activation matching. In *ECCV*, 2018. 3
- [29] Stephen James, Andrew J. Davison, and Edward Johns. Transferring end-to-end visuomotor control from simulation to real world for a multi-stage task. In *CORL*, 2017. 3
- [30] Gregory Kahn, Adam Villafior, Bosen Ding, Pieter Abbeel, and Sergey Levine. Self-supervised deep reinforcement learning with generalized computation graphs for robot navigation. In *ICRA*, 2018. 3
- [31] Liyiming Ke, Xiujun Li, Yonatan Bisk, Ari Holtzman, Zhe Gan, Jingjing Liu, Jianfeng Gao, Yejin Choi, and Siddhartha Srinivasa. Tactical rewind: Self-correction via backtracking in vision-and-language navigation. In *CVPR*, 2019. 3

- [32] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. AI2-THOR: An Interactive 3D Environment for Visual AI. *arXiv*, 2017. 2, 3
- [33] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 1
- [34] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context. In *ECCV*, 2014. 2
- [35] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I. Jordan. Learning transferable features with deep adaptation networks. In *ICML*, 2015. 3
- [36] Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zsolt Kira, Richard Socher, and Caiming Xiong. Self-monitoring navigation agent via auxiliary progress estimation. In *ICLR*, 2019. 3
- [37] Chih-Yao Ma, Zuxuan Wu, Ghassan AlRegib, Caiming Xiong, and Zsolt Kira. The regretful agent: Heuristic-aided navigation through progress estimation. In *CVPR*, 2019. 3
- [38] Piotr Mirowski, Matthew Koichi Grimes, Mateusz Malinowski, Karl Moritz Hermann, Keith Anderson, Denis Teplyashin, Karen Simonyan, Koray Kavukcuoglu, Andrew Zisserman, and Raia Hadsell. Learning to navigate in cities without a map. In *NeurIPS*, 2018. 3
- [39] Piotr Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andrew J. Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray Kavukcuoglu, Dharshan Kumaran, and Raia Hadsell. Learning to navigate in complex environments. In *ICLR*, 2017. 3
- [40] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *ICML*, 2016. 6
- [41] Arsalan Mousavian, Alexander Toshev, Marek Fiser, Jana Kosecka, and James Davidson. Visual representations for semantic target driven navigation. In *ECCV Workshop on Visual Learning and Embodied Agents in Simulation Environments*, 2018. 3
- [42] Matthias Mueller, Alexey Dosovitskiy, Bernard Ghanem, and Vladlen Koltun. Driving policy transfer via modularity and abstraction. In *CORL*, 2018. 3
- [43] Adithyavairavan Murali, Tao Chen, Kalyan Vasudev Alwala, Dhiraj Gandhi, Lerrel Pinto, Saurabh Gupta, and Abhinav Gupta. Pyrobot: An open-source robotics framework for research and benchmarking. *arXiv*, 2019. 3, 5
- [44] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *ICRA*, 2018. 3
- [45] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke S. Zettlemoyer. Deep contextualized word representations. In *NAACL*, 2018. 1
- [46] Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. Virtualhome: Simulating household activities via programs. In *CVPR*, 2018. 3
- [47] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+ questions for machine comprehension of text. In *EMNLP*, 2016. 2
- [48] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 1, 6
- [49] Charles Richter and Nicholas Roy. Safe visual navigation via deep learning and novelty detection. In *RSS*, 2017. 3
- [50] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *IJCV*, 2014. 2
- [51] Fereshteh Sadeghi and Sergey Levine. CAD2RL: real single-image flight without a single real image. In *RSS*, 2017. 3
- [52] Nikolay Savinov, Alexey Dosovitskiy, and Vladlen Koltun. Semi-parametric topological memory for navigation. In *ICLR*, 2018. 3
- [53] Manolis Savva, Angel X. Chang, Alexey Dosovitskiy, Thomas A. Funkhouser, and Vladlen Koltun. MINOS: multimodal indoor simulator for navigation in complex environments. *arXiv*, 2017. 2, 3
- [54] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A platform for embodied ai research. In *ICCV*, 2019. 2, 3
- [55] Shital Shah, Debadeepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In *Field and Service Robotics*, 2017. 3
- [56] Shaojie Shen, Nathan Michael, and Vijay Kumar. Autonomous multi-floor indoor navigation with a computationally constrained mav. In *ICRA*, 2011. 3
- [57] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Josh Susskind, Wenda Wang, and Russell Webb. Learning from simulated and unsupervised images through adversarial training. In *CVPR*, 2017. 3
- [58] Siddhartha S. Srinivasa, Patrick Lancaster, Johan Michalove, Matt Schmittle, Colin Summers, Matthew Rickett, Joshua R. Smith, Sanjiban Choudhury, Christoforos Mavrogiannis, and Fereshteh Sadeghi. MuSHR: A low-cost, open-source robotic racecar for education and research. *arXiv*, 2019. 3
- [59] Baochen Sun, Jiashi Feng, and Kate Saenko. Return of frustratingly easy domain adaptation. In *AAAI*, 2016. 3
- [60] Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. In *ICLR*, 2017. 3
- [61] Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and Vincent Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots. In *RSS*, 2018. 3
- [62] Charles Thorpe, Martial H. Hebert, Takeo Kanade, and Steven A. Shafer. Vision and navigation for the carnegiemellon navlab. *TPAMI*, 1988. 3
- [63] Joshua Tobin, Rachel H Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *IROS*, 2017. 3

- [64] Chris Urmson, Joshua Anhalt, Drew Bagnell, Christopher Baker, Robert Bittner, M. N. Clark, John Dolan, Dave Duggins, Tugrul Galatali, Chris Geyer, Michele Gittleman, Sam Harbaugh, Martial Hebert, Thomas M. Howard, Sascha Kol-ski, Alonzo Kelly, Maxim Likhachev, Matt McNaughton, Nick Miller, Kevin Peterson, Brian Pilnick, Raj Rajkumar, Paul Rybski, Bryan Salesky, Young-Woo Seo, Sanjiv Singh, Jarrod Snider, Anthony Stentz, William “Red” Whittaker, Ziv Wolkowicki, Jason Ziglar, Hong Bae, Thomas Brown, Daniel Demitrish, Bakhtiar Litkouhi, Jim Nickolaou, Var-sha Sadekar, Wende Zhang, Joshua Struble, Michael Taylor, Michael Darms, and Dave Ferguson. Autonomous driving in urban environments: Boss and the urban challenge. *J of Field Robotics*, 2008. 3
- [65] Laurens van der Maaten and Geoffrey E. Hinton. Visualizing data using t-sne. *JMLR*, 2008. 7
- [66] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. Glue: A multi-task benchmark and analysis platform for natural language un-derstanding. In *BlackboxNLP@EMNLP*, 2018. 2
- [67] Xin Wang, Qiuyuan Huang, Asli Celikyilmaz, Jianfeng Gao, Dinghan Shen, Yuan-Fang Wang, William Yang Wang, and Lei Zhang. Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation. In *CVPR*, 2019. 3
- [68] Xin Wang, Wenhan Xiong, Hongmin Wang, and William Yang Wang. Look before you leap: Bridging model-free and model-based reinforcement learning for planned-ahead vision-and-language navigation. In *ECCV*, 2018. 2, 3
- [69] Mitchell Wortsman, Kiana Ehsani, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Learning to learn how to learn: Self-adaptive visual navigation using meta-learning. In *CVPR*, 2019. 2, 3, 5
- [70] Yi Wu, Yuxin Wu, Georgia Gkioxari, and Yuandong Tian. Building generalizable agents with a realistic and rich 3d environment. *arXiv*, 2018. 3
- [71] Yi Wu, Yuxin Wu, Aviv Tamar, Stuart Russell, Georgia Gkioxari, and Yuandong Tian. Bayesian relational memory for semantic visual navigation. In *ICCV*, 2019. 3
- [72] Fei Xia, Amir R. Zamir, Zhiyang He, Alexander Sax, Jiten-dra Malik, and Silvio Savarese. Gibson env: Real-world per-ception for embodied agents. In *CVPR*, 2018. 2, 3
- [73] Huazhe Xu, Yang Gao, Fisher Yu, and Trevor Darrell. End-to-end learning of driving models from large-scale video datasets. In *CVPR*, 2017. 3
- [74] Claudia Yan, Dipendra Kumar Misra, Andrew Bennett, Aaron Walsman, Yonatan Bisk, and Yoav Artzi. CHALET: cornell house agent learning environment. *arXiv*, 2018. 3
- [75] Wei Yang, Xiaolong Wang, Ali Farhadi, Abhinav Gupta, and Roozbeh Mottaghi. Visual semantic navigation using scene priors. In *ICLR*, 2019. 3
- [76] Yang Zhang, Philip David, and Boqing Gong. Curricu-lum domain adaptation for semantic segmentation of urban scenes. In *ICCV*, 2017. 3
- [77] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 3, 8
- [78] Xinge Zhu, Hui Zhou, Ceyuan Yang, Jianping Shi, and Dahua Lin. Penalizing top performers: Conservative loss for semantic segmentation adaptation. In *ECCV*, 2018. 3
- [79] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J. Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *ICRA*, 2017. 3