

# Benchmarking Domain Randomisation for Visual Sim-to-Real Transfer

Raghad Alghonaim<sup>1</sup> and Edward Johns<sup>1</sup>

**Abstract**—Domain randomisation is a very popular method for visual sim-to-real transfer in robotics, due to its simplicity and ability to achieve transfer without any real-world images at all. Nonetheless, a number of design choices must be made to achieve optimal transfer. In this paper, we perform a comprehensive benchmarking study on these different choices, with two key experiments evaluated on a real-world object pose estimation task. First, we study the rendering quality, and find that a small number of high-quality images is superior to a large number of low-quality images. Second, we study the type of randomisation, and find that both distractors and textures are important for generalisation to novel environments.

## I. INTRODUCTION

In recent years, deep learning has been successfully applied to a range of robotics applications, with particular success in those which require visual observations for control [1], [2], [3], [4]. Here, convolutional neural networks (CNNs) enable learning of task-specific visual features directly from data, avoiding the need for laborious task-specific engineering, and potentially outperforming methods using hand-crafted visual representations. However, the reliance of deep learning on large labelled datasets presents a significant challenge.

One of the most promising solutions is sim-to-real transfer, where training is performed in simulation, and a controller is transferred directly to the real world. Of the many sim-to-real transfer methods for vision, domain randomisation [5], [6], [3] is the most popular, since it is simple to implement, and can achieve zero-shot transfer without any real-world data. However, despite its popularity, there are no significant works which benchmark the different types of domain randomisation for visual sim-to-real transfer. Recent benchmarking work has studied sim-to-real for dynamics [7], but the visual sim-to-real problem is a distinct problem and is typically treated independently from dynamics.

In this paper, we study a range of different design choices and empirically evaluate their effect on a real-world task. We divide the paper into two distinct experiments, which study two modes of design choices. In the first experiment, we consider the case when the scene content is known in advance, and the primary role of sim-to-real is to model the effect of illumination and image noise on the observed image. Here, we study how the quality of the rendered images, in terms of the fidelity of graphics pipeline, affects the sim-to-real performance. We also study the optimal trade-off between low-quality and high-quality images, given a fixed

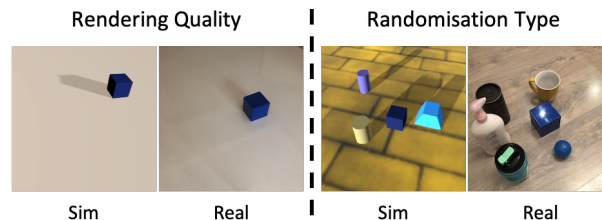


Fig. 1. An overall illustration of our experiments. Left: we study the effect of the rendering quality on the sim-to-real transfer performance with known environment settings. Right: we examine the impact of different randomisation types on the transferability of models trained in simulation and tested in challenging real-world scenarios.

amount of rendering time. In the second experiment, we consider the case when the scene content is not known in advance, and the role of sim-to-real is to achieve robustness to the unknown, such as illumination conditions and clutter. Here, we study how different types of randomisation, such as colours and textures, affect the sim-to-real performance.

We evaluated both experiments with a 6D object pose estimation task, with a manually-labelled real-world dataset. This not only evaluates sim-to-real for a pose estimation module within a wider control pipeline, but also is a proxy to end-to-end control methods [3] which implicitly localise important objects within an image. Our results highlight the importance of high-quality images for sim-to-real transfer, and the importance of randomising both distractors and textures. Our supplementary video gives an overview of our experimental setup, as well as visualisations of results\*.

## II. RELATED WORK

Deep neural networks are known to be sample-inefficient, as a massive amount of annotated images is required for the training process to succeed. Due to the expensive data collection on real robots [8], researchers have shifted towards using simulators to generate the training images. Nonetheless, models trained with synthetic images only are not capable of directly transferring to the real world [9], because of their various physical and visual differences, a problem referred to as the reality gap. Sim-to-real transfer is an active area of research to bridge the gap with two main approaches, *Domain Adaptation* and *Domain Randomisation*. To close the reality gap, sim-to-real transfer methods are extended to both visual [3], [5], [6] and dynamics [7], [10], [11], [12], [13] solutions. The remainder of this section focuses on reviewing visual sim-to-real approaches.

<sup>1</sup> The Robot Learning Lab at Imperial College London {raa318, e.johns}@imperial.ac.uk

\*[www.robot-learning.uk/benchmarking-domain-randomisation](http://www.robot-learning.uk/benchmarking-domain-randomisation)

### A. Domain Adaptation

Domain adaptation is an approach to overcome the reality gap by learning transferable representations across different domains, either by minimising a distance measure between the two distributions [14], [15], [16] or by learning pixel-level representations to make the simulated images as close as possible to their real-world counterparts [17], [18]. While domain adaptation showed promising results in the literature, incorporating real-world images is necessary for successful transfers. In this work, we focus on the zero-shot sim-to-real transfer solutions that do not require finetuning with real-world images.

### B. Domain Randomisation

One of the most common forms of sim-to-real transfer is domain randomisation, a simple-yet-promising technique to bridge the reality gap. Instead of collecting the training data from a single simulated environment, the model is exposed to a variety of random environments to enforce domain invariance. One of the earliest successes of domain randomisation is the pioneering work presented in [6], where the authors leveraged the technique to train for a collision-free indoor flight. They showed that a policy trained solely in simulation is capable of adapting to real-world scenarios if exposed to a sufficient amount of randomness at training time. Several other works have then extended on this success for a variety of robotics manipulation tasks [5], [3], [19], [20], [21], [22], [4], [23], [24].

In [3], the researchers investigated the use of domain randomisation with a task that requires hand-eye coordination (robotic grasping task). Their end-to-end approach succeeded in transferring to clear environments but failed to generalise to more cluttered ones. In [9], the focus was on training object detector using domain randomisation. Unlike other researches, they introduced the use of a non-realistic noise in the randomised scenes. Although their zero-shot model produced adequate results, they needed to incorporate real-world data for better performance. The authors in [19] exploited the randomisation of both physical and visual properties to train a dexterous robotic hand to perform an in-hand manipulation, where they further improved their work in [25] by proposing an iterative approach to learn the randomisation parameters' distributions, the so-called Automatic Domain Randomisation [26], [27], [28].

The authors in [5] employed domain randomisation for the task of 2D object detection. Their model achieved a pose estimation error of 1.5 cm upon tested in the real world. Their work, however, was tailored to regressing to the object position relative to the world frame, but not its orientation, which is crucial for the majority of robotics manipulation tasks.

In general, the main goal of domain randomisation is to find what type of randomisation is effective in different scenarios. Both [3] and [5] found that models are sensitive to the presence of distractors and textures in the simulated scenes. Several other factors have also been investigated in the literature, such as object textures [3], [5], [9], [29], [30],

addition of distractors [3], [5], [9], [29], lights properties [9], [29], target object shapes [31], [29], [30], and sample size [5], [3], [30]. However, none of these works have benchmarked the different choices that must be made to achieve the optimal transfer. In this paper, we empirically evaluate two modes of design choices on a real-world 6D pose estimation task, which serves as a reasonable proxy for evaluating end-to-end control.

## III. METHOD

This work is composed of two distinct experiments that are tailored to different objectives. In this section, we describe the setting shared between both experiments.

### A. Problem Definition

The goal of the experiments is to train two models,  $p_O(I_C)$  and  $q_O(I_C)$ , to map an RGB image  $I$  captured from a camera  $C$  to the 3D position  $(x_O^p, y_O^p, z_O^p)$  of the target object  $O$  and its 4D quaternion  $(x_O^q, y_O^q, z_O^q, w_O^q)$  relative to the camera frame. We chose to regress to quaternions due to their compact form, although there are other possible representations for 3D rotations [32]. While we could train a single network to regress to the full 6D pose of the target object, as is often done in practice [33], we trained two separate networks to avoid the laborious process of finding the optimal balancing between the position and the orientation terms in the loss function, which would distract us from the main point of the paper. The training images are rendered in simulation, and depending on the experiment, they can be of different quality levels, where distractors might also be present in the environment. At test time, the trained model is used without any additional finetuning with real-world images.

### B. Training (Synthetic) Data Collection

We used Unity3D [34] to render all of our training images. We chose to work with Unity as it provides several render pipelines that are used to generate images of different quality levels. Within each sample scene, we define a plane upon which the target object is positioned. Further, we have a camera that moves and rotates in all the three dimensions to capture the target object in different poses, along with some directional lights to illuminate the scene. Depending on the experiment, we randomise different parameters while collecting the data to provide enough variability in the training dataset. For all experiments, however, the following aspects are randomised for each generated example:

- The colour of the target object and its background
- The position, orientation, field of view, and focal length of the camera
- The number of lights, their positions, orientations, and specular characteristics

Table I shows the full list of the randomisation parameters along with their values. The additional experiment-specific parameters are discussed in their respective sections separately.

TABLE I  
RANDOMISATION PARAMETERS VALUES.

Parameter	Range of values
Camera position	Placed within a box of size $(70 \times 40 \times 50)$ cm around the target object
Camera rotation	$x$ : $[0^\circ, 67.5^\circ]$ $y$ : $[-25^\circ, 75^\circ]$ $z$ : $[-10^\circ, 55^\circ]$
Camera field of view	$\pm 5^\circ$ from the estimated real-world equivalent
Target object colour	$\pm 10\%$ from the estimated real-world equivalent
Number of lights	$[0,3]$ directional lights, depending on the quality level
Lights intensity	HDRP: $[0, 2000]$ lux Built-in renderer: $[0,1]$
Lights position	Placed within a box of size $(6 \times 2.8 \times 6)$ m around the target object
Lights rotation	$x$ : $[5^\circ, 30^\circ]$ $y$ : $[-180^\circ, 180^\circ]$ $z$ : $[-90^\circ, 90^\circ]$

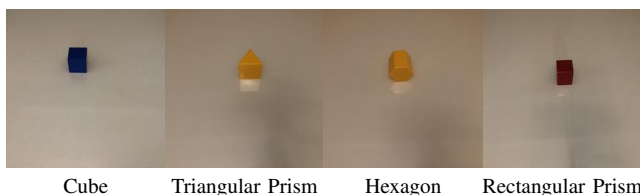


Fig. 2. The four primitive shapes that we used throughout the experiments.

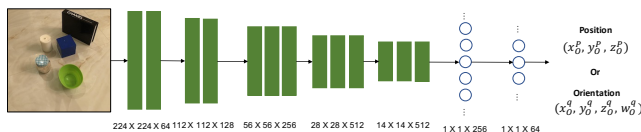


Fig. 3. The end-to-end task of 6D pose estimation. For each target object  $O$ , we train two separate networks,  $p_O(I_C)$  and  $q_O(I_C)$ , that map a single image to the 6D pose of  $O$ . Each input image is downsized to  $(224 \times 224)$  and supplied to two separate networks to be mapped to a tuple of  $(x_O^p, y_O^p, z_O^p)$  that represents the 3D position of the target object and a tuple of  $(x_O^q, y_O^q, z_O^q, w_O^q)$  which represents the target’s 3D orientation.

### C. Testing (Real-world) Data Collection

To address the research questions, we created mesh representations for four primitive shapes, shown in Fig. 2, that we used to evaluate both experiments. To find the accurate ground-truth values, we used *ArUco* [35], an accurate marker-based pose estimation algorithm. For each test example, we placed a marker on top of the target object and carefully validated, with an accurate ruler, the offset from the object’s coordinate frame. The image is then supplied to *ArUco* which outputs the 6D pose value of the object relative to the calibrated camera. Another image is captured for the same scene without the marker to be used as an input to the trained networks.

### D. Network Architecture

For all experiments, we used a modified version of the VGG-16 [36] network architecture, which was proposed in [5]. As shown in Fig. 3, the network consists of five groups of convolutional layers each with a kernel size of  $3 \times 3$  and a stride of 1, where dimensionality reduction is performed

after each group. To guarantee a fair comparison, we fixed the hyper-parameters, the seed, and the loss function for all experiments. All models are trained using *Adam* optimiser until convergence with a starting learning rate of  $1e-4$  and a scheduler to decrease it by a factor of 0.1 every 30 epochs. The networks weights are initialised randomly, and both the inputs and outputs are normalised before starting the training process.

### E. Evaluation Metrics

To evaluate the performance of our models, we test them using the collected real-world images and compare the actual labels to the results obtained from the networks. The position error is computed using the Root Mean Squared Error (or RMSE) formula. For the orientation, however, we use Equation 1 [37] to find the angle required to rotate from one quaternion (the label,  $q_i$ ) to the other (the prediction,  $\hat{q}_i$ ).

$$\theta_i = 2 \times \cos^{-1}(|\langle \hat{q}_i, q_i \rangle|) \quad (1)$$

Where  $\langle \hat{q}_i, q_i \rangle$  is the inner product between the label and the prediction of the  $i^{th}$  sample image.

## IV. EXPERIMENT 1: RENDERING QUALITY

This experiment aims at studying the impact of the simulator’s fidelity on the overall model’s performance in the real world. More precisely, our focus is on finding answers to the following three questions:

- 1) How critical is the quality of the simulator for achieving successful sim-to-real transfers?
- 2) What is the effect that each simulation parameter has on the overall sim-to-real transfer performance?
- 3) What is the optimal trade-off of low-quality and high-quality images given a fixed amount of rendering time?

### A. Experimental Setup

To achieve this, we use two renderers provided by Unity: the Built-in and the High Definition Render Pipeline (HDRP hereafter) [38] to render images of eight quality levels. As presented in Table II, these levels vary in terms of the following aspects:

- **Directional lights:** we enable them for all quality levels except the baseline (level 1).
- **Shadows:** we alter the quality of the shadows for each simulation level to range from low to high quality.
- **Anti-aliasing:** anti-aliasing is a technique used to overcome the staircase effect that usually occurs in physical simulators by smoothing the object boundaries and corners.
- **Dithering effect:** dithering is a form of noise that is intentionally applied to rendered images to eliminate unwanted graphical issues, such as colour banding.
- **Render pipeline:** the renderer used to generate the dataset.

In Fig. 4, we visually illustrate the differences between the eight simulation quality levels. In this experiment, we assumed that the colours of both the target object and its

TABLE II  
AN OVERALL COMPARISON BETWEEN ALL THE CONSIDERED LEVELS OF SIMULATION.

Quality level \ Feature	Directional Lights	Shadows	Anti-aliasing	Dithering effect	Render Pipeline
1	No	No	No	No	Built-In
2	Yes	No	No	No	
3	Yes	Low-quality	No	No	
4	Yes	Medium-quality	No	No	HDRP
5	Yes	Medium-quality	Low-quality	No	
6	Yes	High-quality	Low-quality	No	
7	Yes	High-quality	High-quality	No	
8	Yes	High-quality	High-quality	Yes	

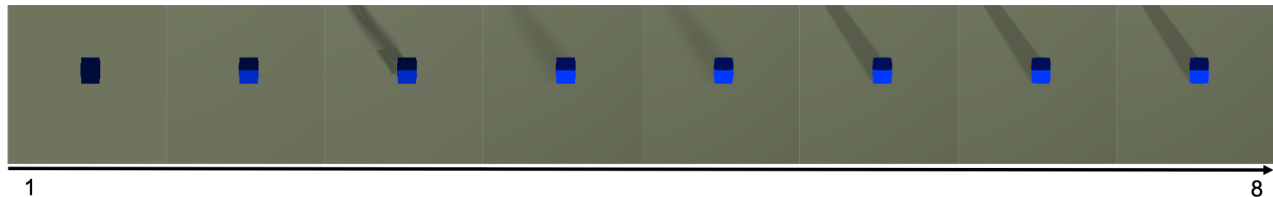


Fig. 4. Sample images from each quality level. From left to right: low-quality (level 1) to high-quality (level 8).

TABLE III  
THE AVERAGE ERROR ACROSS FOUR TARGET OBJECTS.

Quality level	$xy$ -position (cm)	$z$ -position (cm)	Orientation (degrees)
1	$4.37 \pm 4.45$	$12.46 \pm 8.06$	$13.0 \pm 4.72$
2	$1.44 \pm 1.68$	$3.44 \pm 2.87$	$4.57 \pm 3.38$
3	$1.40 \pm 1.53$	$3.22 \pm 3.05$	$4.14 \pm 2.97$
4	$0.83 \pm 1.10$	$2.73 \pm 3.60$	$3.52 \pm 2.26$
5	$0.70 \pm 0.62$	$2.67 \pm 2.22$	$3.34 \pm 2.41$
6	$0.63 \pm 0.63$	$2.50 \pm 2.14$	$3.26 \pm 2.40$
7	$0.61 \pm 0.62$	$2.43 \pm 2.07$	$3.20 \pm 2.57$
8	<b><math>0.53 \pm 0.52</math></b>	<b><math>2.38 \pm 1.93</math></b>	<b><math>3.08 \pm 2.15</math></b>

background are known at training time, and no distractors are present in the environment. For each object, we collected eight training datasets for the different quality levels we consider, each with 30,000 randomised synthetic images. We tested the models with 100 real-world testing images, per object, that are collected with a variety of poses and lighting conditions.

### B. Results

In Table III, we present the average error achieved by each quality level across the four target objects. We reported the  $xy$  and the  $z$  position errors separately to study the effect on the different tasks independently:  $xy$  position (across the image plane),  $z$  position (perpendicular to the image plane), and 3D orientation.

**Altering the Rendering Quality** We can conclude from Table III that, as we increase the quality of the training data, the models performed consistently better, which answers our first question.

**Directional Lights Effect** As we can deduce from the first two rows of Table III, we were able to reduce the error

by  $\approx 68.4\%$  after enabling directional lights, which proves their importance for achieving successful transfer.

**Shadows Effect** Without shadows, as demonstrated in the second and third rows of Table III, the models performed the worst when tested in the real world. Moreover, as we increased the quality of the shadows (level 3 vs. level 4), the performance of the models has significantly improved, which indicates the importance of the quality of the shadows.

**Anti-aliasing Effect** Recall from Table II that the only difference between level 4 and level 5 images is the use of anti-aliasing. As we can conclude from Table III, models trained with level 5 images succeeded in outperforming the ones trained with no anti-aliasing enabled (i.e. level 4), which shows the positive impact of anti-aliasing on the overall models' performance. However, since the anti-aliasing used with level 5 images is of low-quality, it might fail, for some objects with complex shapes, to fix the staircase issues.

**Dithering Effect** The last two rows of Table III demonstrate the results before and after enabling this effect, where we can see that dithering plays a vital role in boosting the transfer performance.

### C. Combining High-quality and Low-quality Images

In the previous set of experiments, we showed that models trained with high-quality images have superior performance compared to the ones trained with lower-quality levels. Nonetheless, rendering high-quality images comes with an additional computational cost. We found with experiments that rendering one image from level 1 takes an average of only 8.4 ms, while generating one high-quality image (level 8) requires more than triple the time, with an average of 27.2 ms. Rather than training with images of only one level of quality, we could also train with both high-quality and low-quality images, to combine the merits of both fast and

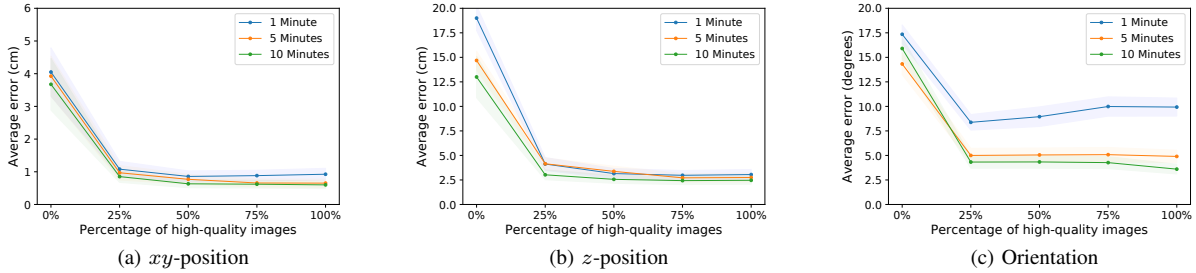


Fig. 5. The results obtained after training the models with combinations of high-quality and low-quality training images, under three different data collection time constraints.

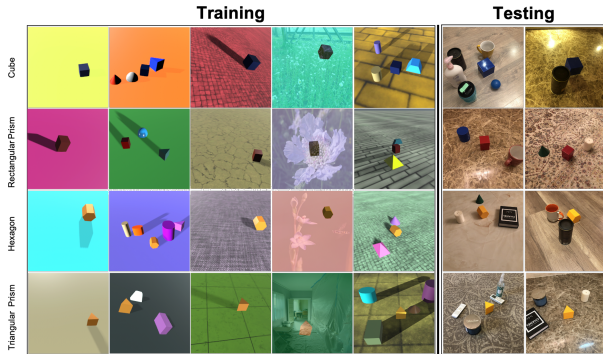


Fig. 6. Left: high-quality synthetic images collected with random colours, textures and distractors. Right: real-world testing images labelled using ArUco.

accurate rendering.

Fig. 5 shows the results obtained after training several models under different time constraints. For each experiment, we fixed the data collection time and varied the percentage of each type of images. As we can deduce from the graphs, the overall real-world test error dropped significantly after incorporating as minimum as 25% of high-quality images. Surprisingly, however, as we increased the percentage of high-quality images, the performance of the models exhibited no clear trend, as it sometimes stabilised, reduced, or improved slightly. The overall conclusion is therefore that if a significant number of images can be collected (i.e. no time constraint), there is no harm from rendering all images with high quality. Otherwise, if only a small number of images can be collected, including some high-quality ones (e.g. 25%) should significantly improve the results.

## V. EXPERIMENT 2: RANDOMISATION TYPE

In this experiment, we examine the performance of models trained with high-quality synthetic images, when tested in varied and cluttered real-world environments. More precisely, we strive to assess the significance of different randomisation settings to the models’ transferability to the real world. Similar to the previous experiment, the focus here is on evaluating the sim-to-real performance, rather than studying the different settings of 6D pose estimation.

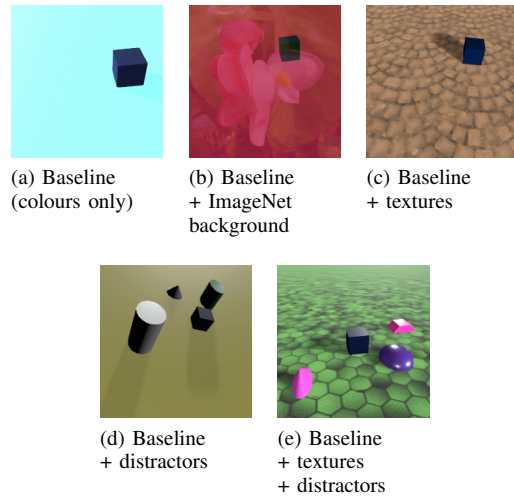


Fig. 7. A sample from each of the five datasets collected to study the impact of the different randomisation settings on the sim-to-real transfer performance. In these examples, the target object is the blue cube.

### A. Experimental Setup

Unlike the previous experiment, in this, we assumed that the models have no prior knowledge about the background of the target object (i.e. it can be of any texture or colour). In addition, distractors of any shape, colour, and size are presented at test time.

A summary of the approach is provided in Fig. 6. The models were trained using high-quality simulated images (left) and evaluated in the real-world (right). Random geometrical shapes were used as distractors at training time. However, a variety of never-seen distractors are presented in the real-world evaluation images, e.g. a cup, book, and candle. To assess which types of randomisation are most crucial for achieving successful transfer, we collected five training datasets, shown in Fig. 7, which have the following specifications:

- 1) Baseline: images with random solid colours in the background, but no textures or distractors
- 2) Baseline + random ImageNet backgrounds
- 3) Baseline + random background textures
- 4) Baseline + random distractors
- 5) Baseline + random textures and distractors

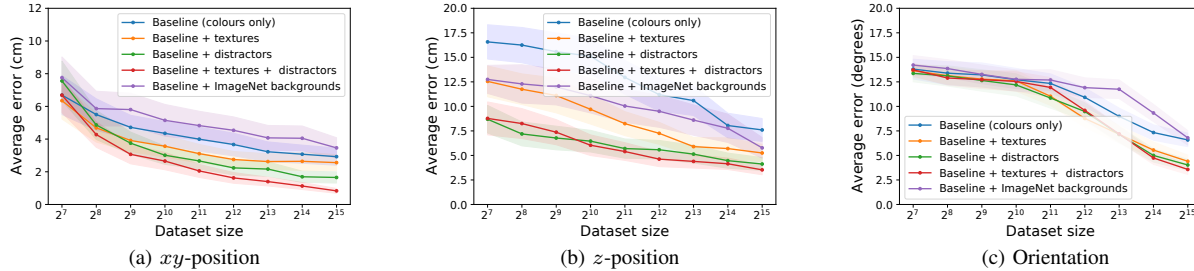


Fig. 8. The sensitivity of the models to distractors, textures, and number of images. The results are averaged across the four target objects.

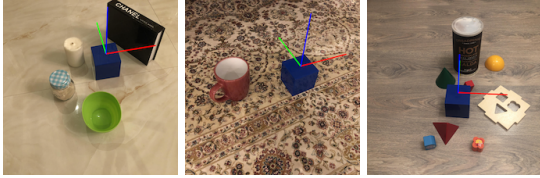


Fig. 9. Examples where the model succeeds in accurately estimating the 6D pose of the target object (cube). The axes represent the cube’s local frame, and their colours correspondences are:  $x$ : red,  $y$ : green, and  $z$ : blue.

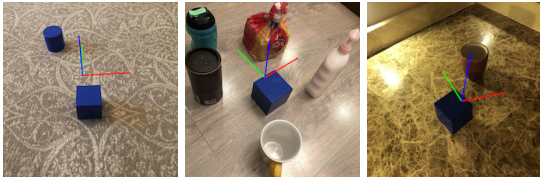


Fig. 10. Examples where the model fails to locate the target. The model missed the object when a distractor of the same colour is presented in the environment and with extremely cluttered or dark scenes. The axes are the cube’s local frame ( $x$ : red,  $y$ : green,  $z$ : blue).

The colour of the target object was uniformly sampled around our best estimate of its colour in the real world. All the five datasets were collected with high-quality anti-aliasing, high-quality shadows, and dithering effect. For each target object, we collected 110 real-world testing images with a variety of backgrounds and distractors. We did not control for lighting and shadows in the real world, meaning that some images were captured from relatively dark scenes. For this experiment, we focused on answering the following three questions:

- 1) What elements are most important to randomise while collecting the training data?
- 2) How does the performance of the models change as we increase the training dataset size?
- 3) How robust are the trained models to never-seen backgrounds and distractors?

## B. Results

**Type of Randomisation** We assessed the sensitivity of our models to both textures and distractors. We found, as shown in Fig. 8, that the models are sensitive to both factors for the three tasks. However, we noticed that for position

estimation, training with distractors in the environment (the green line in Fig. 8) is of more importance than having textures, which is concluded from the fact that models trained with textures and no distractors had larger errors compared to the ones trained with distractors only. In general, both the baseline and the ImageNet datasets failed with the worst performance compared to the rest of the models. In contrast, the full dataset attained the best performance in all cases.

**Varying the Dataset Size** To answer the second question, we trained several models on different dataset sizes, ranging from  $2^7$  to  $2^{15}$  images. The graphs in Fig. 8 show that as we increase the training dataset size, the average real-world error gradually decreased for all the five datasets, as would be expected.

**Qualitative Results** Overall, we observed that our models, when trained with both textures and distractors, are robust to changing environments and can generalise to cases with never-seen distractors and backgrounds. In Fig. 9, we show some examples when the model precisely estimated the 6D pose of the target object, in spite of the novel backgrounds and distractors. Despite the promising results, the model was not able to accurately estimate the pose of the target object, as shown in Fig. 10, when another distractor of the same colour is presented in the environment. Further, extremely cluttered environments and dark scenes posed a challenge to the model.

## VI. CONCLUSION

This paper investigated two different design choices in domain randomisation for visual sim-to-real transfer: the quality of the renderer, and the type of randomisation. For the first set of experiments, we found there to be a very strong relationship between the quality of the renderer, and the sim-to-real performance. We also proposed to combine low-quality images with high-quality images, and found that for the same overall rendering time, it is more important to have a high percentage of high-quality images, than low-quality images. For the second set of experiments, we showed that randomising both distractor objects and background textures are important for generalising to novel environments. These conclusions can now be used by others in designing their own domain randomisation datasets, with a view towards achieving optimal sim-to-real performance for a given amount of dataset generation time.

## REFERENCES

- [1] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *The Journal of Machine Learning Research*, 2016.
- [2] E. Johns, "Coarse-to-fine imitation learning: Robot manipulation from a single demonstration," in *International Conference on Robotics and Automation (ICRA)*, 2021.
- [3] S. James, A. J. Davison, and E. Johns, "Transferring end-to-end visuomotor control from simulation to real world for a multi-stage task," in *Conference on Robot Learning (CoRL)*, 2017.
- [4] K. Bousmalis, A. Irpan, P. Wohlhart, Y. Bai, M. Kelcey, M. Kalakrishnan, L. Downs, J. Ibarz, P. Pastor, K. Konolige, S. Levine, and V. Vanhoucke, "Using simulation and domain adaptation to improve efficiency of deep robotic grasping," *CoRR*, vol. abs/1709.07857, 2017. [Online]. Available: <http://arxiv.org/abs/1709.07857>
- [5] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2017, pp. 23–30.
- [6] F. Sadeghi and S. Levine, "Cad2rl: Real single-image flight without a single real image," *arXiv preprint arXiv:1611.04201*, 2016.
- [7] E. Valassakis, Z. Ding, and E. Johns, "Crossing the gap: A deep dive into zero-shot sim-to-real transfer for dynamics," in *International Conference on Intelligent Robots and Systems (IROS)*, 2020.
- [8] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen, "Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection," *The International Journal of Robotics Research*, vol. 37, no. 4-5, pp. 421–436, 2018.
- [9] J. Tremblay, A. Prakash, D. Acuna, M. Brophy, V. Jampani, C. Anil, T. To, E. Cameracci, S. Boochoon, and S. Birchfield, "Training deep networks with synthetic data: Bridging the reality gap by domain randomization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 969–977.
- [10] Y.-Y. Tsai, H. Xu, Z. Ding, C. Zhang, E. Johns, and B. Huang, "Droid: Minimizing the reality gap using single-shot human demonstration," *IEEE Robotics and Automation Letters*, 2021.
- [11] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," *CoRR*, vol. abs/1710.06537, 2017. [Online]. Available: <http://arxiv.org/abs/1710.06537>
- [12] A. Rajeswaran, S. Ghotra, S. Levine, and B. Ravindran, "Epopt: Learning robust neural network policies using model ensembles," *CoRR*, vol. abs/1610.01283, 2016. [Online]. Available: <http://arxiv.org/abs/1610.01283>
- [13] W. Yu, C. K. Liu, and G. Turk, "Preparing for the unknown: Learning a universal policy with online system identification," *CoRR*, vol. abs/1702.02453, 2017. [Online]. Available: <http://arxiv.org/abs/1702.02453>
- [14] M. Long, Y. Cao, J. Wang, and M. I. Jordan, "Learning transferable features with deep adaptation networks," *arXiv preprint arXiv:1502.02791*, 2015.
- [15] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7167–7176.
- [16] B. Sun and K. Saenko, "Deep coral: Correlation alignment for deep domain adaptation," in *European conference on computer vision*. Springer, 2016, pp. 443–450.
- [17] K. Bousmalis, A. Irpan, P. Wohlhart, Y. Bai, M. Kelcey, M. Kalakrishnan, L. Downs, J. Ibarz, P. Pastor, K. Konolige, S. Levine, and V. Vanhoucke, "Using simulation and domain adaptation to improve efficiency of deep robotic grasping," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 4243–4250.
- [18] S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell, and K. Bousmalis, "Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 627–12 637.
- [19] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al., "Learning dexterous in-hand manipulation," *The International Journal of Robotics Research*, vol. 39, no. 1, pp. 3–20, 2020.
- [20] F. Sadeghi, A. Toshev, E. Jang, and S. Levine, "Sim2real view invariant visual servoing by recurrent control," *arXiv preprint arXiv:1712.07642*, 2017.
- [21] J. Matas, S. James, and A. J. Davison, "Sim-to-real reinforcement learning for deformable object manipulation," *CoRR*, vol. abs/1806.07851, 2018. [Online]. Available: <http://arxiv.org/abs/1806.07851>
- [22] S. James, M. Bloesch, and A. J. Davison, "Task-embedded control networks for few-shot imitation learning," in *Conference on Robot Learning*. PMLR, 2018, pp. 783–795.
- [23] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel, "Asymmetric actor critic for image-based robot learning," *CoRR*, vol. abs/1710.06542, 2017. [Online]. Available: <http://arxiv.org/abs/1710.06542>
- [24] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. D. Ratliff, and D. Fox, "Closing the sim-to-real loop: Adapting simulation randomization with real world experience," *CoRR*, vol. abs/1810.05687, 2018. [Online]. Available: <http://arxiv.org/abs/1810.05687>
- [25] I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, et al., "Solving rubik's cube with a robot hand," *arXiv preprint arXiv:1910.07113*, 2019.
- [26] C. Heindl, S. Zambal, and J. Scharinger, "Learning to predict robot keypoints using artificially generated images," in *2019 24th IEEE International Conference on Emerging Technologies and Factory Automation (ETFA)*, 2019, pp. 1536–1539.
- [27] N. Ruiz, S. Schuler, and M. Chandraker, "Learning to simulate," *arXiv preprint arXiv:1810.02513*, 2018.
- [28] C. Heindl, L. Brunner, S. Zambal, and J. Scharinger, "Blendtorch: A real-time, adaptive domain randomization library," *arXiv preprint arXiv:2010.11696*, 2020.
- [29] 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*, vol. 36, no. 1, p. 1–14, Feb 2020. [Online]. Available: <http://dx.doi.org/10.1109/TRO.2019.2942989>
- [30] A. Dehban, J. Borrego, R. Figueiredo, P. Moreno, A. Bernardino, and J. Santos-Victor, "The impact of domain randomization on object detection: A case study on parametric shapes and synthetic textures\*," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 2593–2600.
- [31] J. Tobin, L. Biewald, R. Duan, M. Andrychowicz, A. Handa, V. Kumar, B. McGrew, J. Schneider, P. Welinder, W. Zaremba, and P. Abbeel, "Domain randomization and generative models for robotic grasping," 2018.
- [32] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li, "On the continuity of rotation representations in neural networks," 2020.
- [33] A. Kendall, M. Grimes, and R. Cipolla, "Convolutional networks for real-time 6-dof camera relocalization," *CoRR*, vol. abs/1505.07427, 2015. [Online]. Available: <http://arxiv.org/abs/1505.07427>
- [34] U. Technologies, "Unity3d." [Online]. Available: <https://unity3d.com/>
- [35] F. J. Romero-Ramirez, R. Muñoz-Salinas, and R. Medina-Carnicer, "Speeded up detection of squared fiducial markers," *Image and vision Computing*, vol. 76, pp. 38–47, 2018.
- [36] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014.
- [37] S. Mahendran, H. Ali, and R. Vidal, "3d pose regression using convolutional neural networks," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 2174–2182.
- [38] Unity3D, "High definition render pipeline/built-in render pipeline comparison." [Online]. Available: <https://docs.unity3d.com/Packages/com.unity.render-pipelines.high-definition@7.1/manual/Feature-Comparison.html>