# *3DB*: A Framework for Debugging Computer Vision Models

Guillaume Leclerc[†]
LECLERC@MIT.EDU
MIT [*]

Hadi Salman[†]
HADY@MIT.EDU
MIT [*]

Andrew Ilyas[†]
AILYAS@MIT.EDU
MIT

Sai Vemprala
SAIHV@MICROSOFT.COM
Microsoft Research

Logan Engstrom
ENGSTROM@MIT.EDU
MIT

Vibhav Vineet
VIVINEET@MICROSOFT.COM
Microsoft Research

Kai Xiao
KAIX@MIT.EDU
MIT

Pengchuan Zhang
PENZHAN@MICROSOFT.COM
Microsoft Research

Shibani Santurkar
SHIBANI@MIT.EDU
MIT

Greg Yang
GE.YANG@MICROSOFT.COM
Microsoft Research

Ashish Kapoor
AKAPOOR@MICROSOFT.COM
Microsoft Research

Aleksander Mądry
MADRY@MIT.EDU
MIT

## Abstract

We introduce *3DB*: an extendable, unified framework for testing and debugging vision models using photorealistic simulation. We demonstrate, through a wide range of use cases, that *3DB* allows users to discover vulnerabilities in computer vision systems and gain insights into how models make decisions. *3DB* captures and generalizes many robustness analyses from prior work, and enables one to study their interplay. Finally, we find that the insights generated by the system transfer to the physical world.

We are releasing *3DB* as a library[1] alongside a set of example analyses[2], guides[3], and documentation[4].

---

# 1   Introduction

Modern machine learning models turn out to be remarkably brittle under distribution shift. Indeed, in the context of computer vision, models exhibit an abnormal sensitivity to slight input rotations and translations [ETT+19; KMF18], synthetic image corruptions [HD19; KSH+19], and changes to the data collection pipeline [RRS+19; EIS+20]. Still, while such brittleness is widespread, it is often hard to understand its root causes, or even to characterize the precise situations in which this unintended behavior arises.

How do we then comprehensively diagnose model failure modes? Stakes are often too high to simply deploy models and collect eventual "real-world" failure cases. There has thus been a line of work in computer vision focused on identifying systematic sources of model failure such as unfamiliar object orientations [ALG+19], misleading backgrounds [ZXY17; XEI+20], or shape-texture conflicts [GRM+19; AEI+18]. These analyses—a selection of which is visualized in Figure 1—reveal patterns or situations that degrade performance of vision models, providing invaluable insights into model robustness. Still, carrying out each such analysis requires its own set of (often complex) tools and techniques, usually accompanied by a significant amount of manual labor (e.g., image editing, style transfer, etc.), expertise, and data cleaning. This prompts the question:

*Can we support reliable discovery of model failures in a systematic, automated, and unified way?*

**Contributions.**    In this work, we propose *3DB*, a framework for automatically identifying and analyzing the failure modes of computer vision models. This framework makes use of a 3D simulator to render realistic scenes that can be fed into any computer vision system. Users can specify a set of transformations to apply to the scene—such as pose changes, background changes, or camera effects—and can also customize and compose them. The system then performs a guided search, evaluation, and aggregation over these user-specified configurations and presents the user with an interactive, user-friendly summary of the model's performance and vulnerabilities. *3DB* is general enough to enable users to, with little-to-no effort, re-discover insights from prior work on robustness to pose, background, and texture bias (cf. Figure 2), among others. Further, while prior studies have largely been focused on examining model sensitivities along a single axis, *3DB* allows users to compose various transformations to understand the interplay between them, while still being able to disentangle their individual effects.
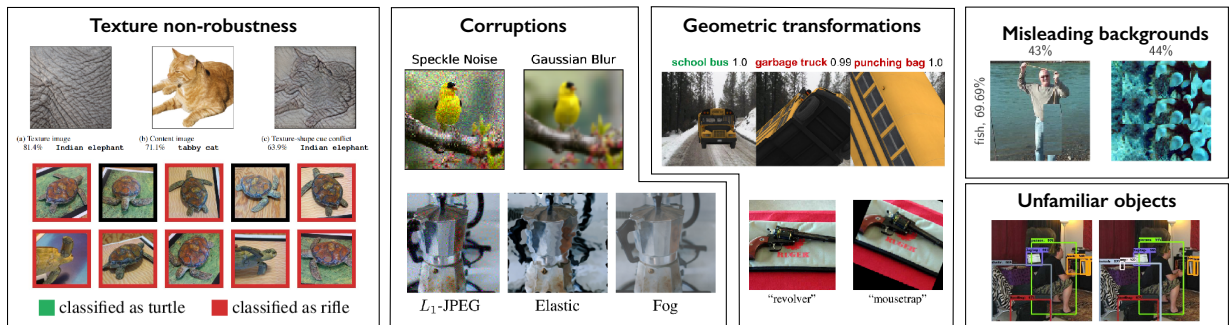


Figure 1: Examples of vulnerabilities of computer vision systems identified through prior in-depth robustness studies. Figures reproduced from [GRM+19; AEI+18; HD19; KSH+19; ALG+19; ETT+19; XEI+20; RZT18].



Figure 2: The *3DB* framework is modular enough to facilitate—among other tasks—efficient rediscovery of all the types of brittleness shown in Figure 1 in an integrated manner. It also allows users to realistically compose transformations (right) while still being able to disentangle the results.

The remainder of this paper is structured into the following three parts: in Section 2 we discuss the design of *3DB*, including the motivating principles, design goals, and concrete architecture used. We highlight how the implementation of *3DB* allows users to quickly experiment, stress-test, and analyze their vision models. Then, in Section 3 we illustrate the utility of *3DB* through a series of case studies uncovering biases in an ImageNet-pretrained classifier. Finally, we show (in Section 4) that the vulnerabilities uncovered with *3DB* correspond to actual failure modes in the physical world (i.e., they are not specific to simulation).

# 2 Designing *3DB*

The goal of *3DB* is to leverage photorealistic simulation in order to effectively diagnose failure modes of computer vision models. To this end, the following set of principles guide the design of *3DB*:

(a) **Generality**: *3DB* should support any type of computer vision model (i.e., not necessarily a neural network) trained on any dataset and task (i.e., not necessarily classification). Furthermore, the framework should support diagnosing non-robustness with respect to any parameterizable three-dimensional scene transformation.

(b) **Compositionality**: Data corruptions and transformations rarely occur in isolation. Thus, *3DB* should allow users to investigate robustness along many different axes simultaneously.

(c) **Physical realism**: The vulnerabilities extracted from *3DB* should correspond to models' behavior in the real (physical) world, and, in particular, not depend on artifacts of the simulation process itself. Specifically, the insights that *3DB* produces should not be affected by a simulation-to-reality gap, and still hold when models are deployed in the wild.

(d) **User-friendliness**: *3DB* should be simple to use and should relay insights to the user in an easy-to-understand manner. Even non-experts should be able to look at the result of a *3DB* experiment and easily understand what the weak points of their model are, as well as gain insight into how the model behaves more generally.

(e) **Scalability**: *3DB* should be performant and parallelizable.

## 2.1 Capabilities and workflow

To achieve the goals articulated above, we design *3DB* in a modular manner, i.e., as a combination of swappable components. This combination allows the user to specify transformations they want to test, search over the space of these transformations, and aggregate the results of this search in a concise way. More specifically, the *3DB* workflow revolves around five steps (visualized in Figure 3):

1. **Setup**: The user collects one or more 3D meshes that correspond to objects the model is trained to recognize, as well as a set of environments to test against.

2. **Search space design**: The user defines a *search space* by specifying a set of transformations (which *3DB* calls *controls*) that they expect the computer vision model to be robust to (e.g., rotations, translations, zoom, etc.). Controls are grouped into "rendered controls" (applied during the rendering process) and "post-processor controls" (applied after the rendering as a 2D image transformation).

3. **Policy-guided search**: After the user has specified a set of controls, *3DB* instantiates and renders a myriad of object configurations derived from compositions of the given transformations. It records the behavior of the ML model on each constructed scene for later analysis. A user-specified *search policy* over the space of all possible combinations of transformations determines the exact scenes for *3DB* to render.

4. **Model loading**: The only remaining step before running a *3DB* analysis is loading the vision model that the user wants to analyze (e.g., a pre-trained classifier or object detection model).

5. **Analysis and insight extraction**: Finally, *3DB* is equipped with a model *dashboard* (cf. Appendix B) that can read the generated log files and produce a user-friendly visualization of the generated insights. By default, the dashboard has three panels. The first of these is failure mode display, which highlights configurations, scenes, and transformations that caused the model to misbehave. The per-object analysis pane allows the user to

inspect the model's performance on a specific 3D mesh (e.g., accuracy, robustness, and vulnerability to groups of transformations). Finally, the aggregate analysis pane extracts insights about the model's performance averaged over all the objects and environments collected and thus allows the user to notice consistent trends and vulnerabilities in their model.

Each of the aforementioned components (the controls, policy, renderer, inference module, and logger) are fully customizable and can be extended or replaced by the user without altering the core code of *3DB*. For example, while *3DB* supports more than 10 types of controls out-of-the-box, users can add custom ones (e.g., geometric transformations) by implementing an abstract function that maps a 3D state and a set of parameters to a new state. Similarly, *3DB* supports debugging classification and object detection models by default, and by implementing a custom evaluator module, users can extend support to a wide variety of other vision tasks and models. We refer the reader to Appendix A for more information on *3DB* design principles, implementation, and scalability.
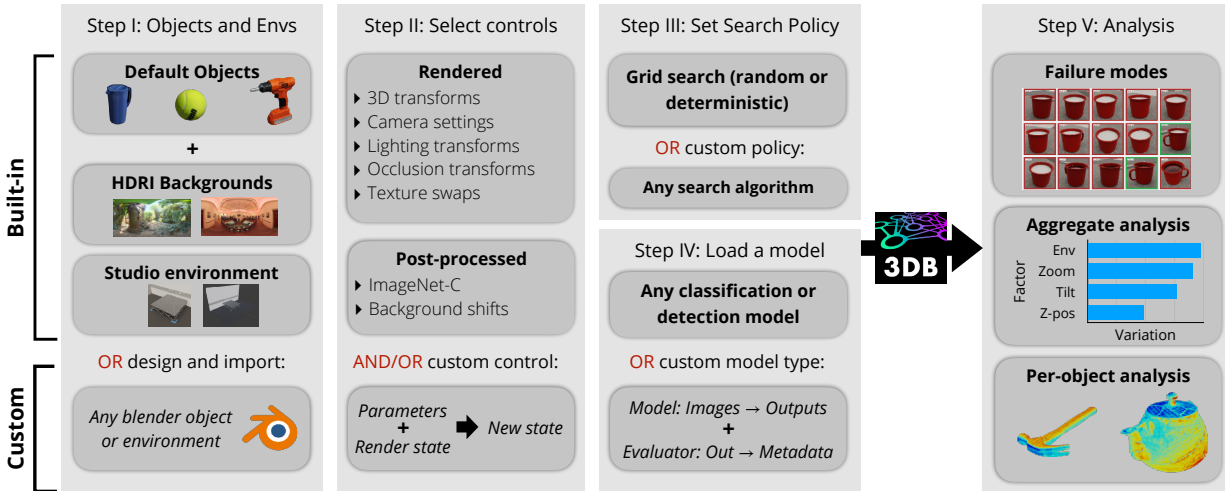


Figure 3: An overview of the *3DB* workflow: First, the user specifies a set of 3D object models and environments to use for debugging. The user also enumerates a set of (in-built or custom) transformations, known as controls, to be applied by *3DB* while rendering the scene. Based on a user-specified search policy over all these controls (and their compositions), *3DB* then selects the exact scenes to render. The computer vision model is finally evaluated on these scenes and the results are logged in a user-friendly manner in a custom dashboard.

# 3 Debugging and Analyzing Models with *3DB*

In this section, we illustrate through case studies how to analyze and debug vision models with *3DB*. In each case, we follow the workflow outlined in Section 2.1—importing the relevant objects, selecting the desired transformations (or constructing custom ones), selecting a search policy, and finally analyzing the results.

In all our experiments, we analyze a ResNet-18 [HZR+15] trained on the ImageNet [RDS+15] classification task (its validation set accuracy is 69.8%). Note that *3DB* is classifier-agnostic (i.e., ResNet-18 can be replaced with any PyTorch classification module), and even supports object detection tasks. For our analysis, we collect 3D models for 16 ImageNet classes (see Appendix E for more details on each experiment). We ensure that in "clean" settings, i.e., when rendered in simple poses on a plain white background, the 3D models are correctly classified at a reasonable rate (cf. Table 1) by our pre-trained ResNet.

## 3.1 Sensitivity to image backgrounds

We begin our exploration by using *3DB* to confirm ImageNet classifiers' reliance on background signal, as pinpointed by several recent in-depth studies [ZML+07; ZXY17; XEI+20]. Out-of-the-box, *3DB* can render 3D models onto HDRI files using image-based lighting; we downloaded 408 such background environments from `hdrihaven.com`. We

|            | banana | baseball | bowl | drill | golf ball | hammer | lemon | mug |
|------------|--------|----------|------|-------|-----------|--------|-------|-----|
| Simulated accuracy (%) | 96.8 | 100.0 | 17.5 | 63.3 | 95.0 | 65.6 | 100.0 | 13.4 |
| ImageNet accuracy (%) | 82.0 | 66.0 | 84.0 | 40.0 | 82.0 | 54.0 | 76.0 | 42.0 |

|            | orange | pitcher base | power drill | sandle | shoe | spatula | teapot | tennis ball |
|------------|--------|--------------|-------------|--------|------|---------|--------|-------------|
| Simulated accuracy (%) | 98.5 | 7.9 | 87.5 | 88.0 | 59.2 | 76.1 | 47.8 | 100.0 |
| ImageNet accuracy (%) | 72.0 | 52.0 | 40.0 | 66.0 | 82.0 | 18.0 | 80.0 | 68.0 |

Table 1: Accuracy of a pre-trained ResNet-18, for each of the 16 ImageNet classes considered, on the corresponding 3D model we collected, rendered at an unchallenging pose on a white background ("Simulated" row); and the subset of the ImageNet validation set corresponding to the class ("ImageNet" row).
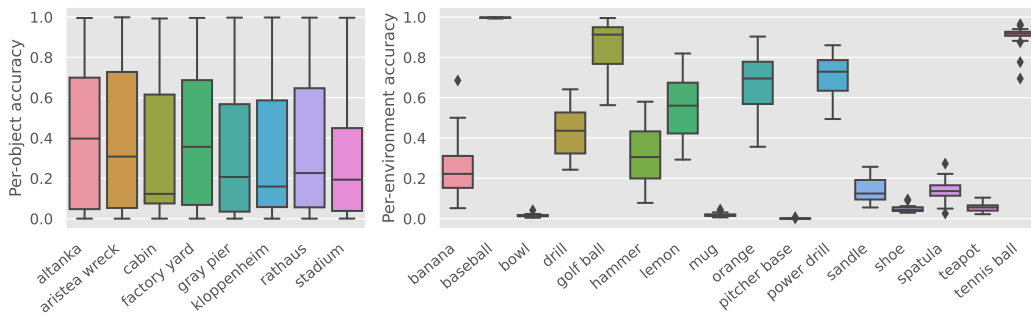


Figure 4: Visualization of accuracy on controls from Section 3.1. **(Left)** We compute the accuracy of the model conditioned on each object-environment pair. For each environment on the x-axis, we plot the variation in accuracy (over the set of possible objects) using a boxplot. We visualize the per-object accuracy spread by including the median line, the first and third quartiles box edges (the interval between which is called the inter-quartile range, IQR), the range, and the outliers (points that are outside the IQR by $3/2|\text{IQR}|$). **(Right)** Using the same format, we track how the classified object (on the x-axis) impacts variation in accuracy (over different environments) on the y-axis.

then used the pre-packaged "camera" and "orientation" controls to render (and evaluate our classifier on) scenes of the pre-collected 3D models at random poses, orientations, and scales on each background. Figure 5 shows some (randomly sampled) example scenes generated by *3DB* for the "coffee mug" model.

**Analyzing a subset of backgrounds.** In Figure 4, we visualize the performance of a ResNet-18 classifier on the 3D models from 16 different ImageNet classes—in random positions, orientations, and scales—rendered onto 20[5] of the collected HDRI backgrounds. One can observe that background dependence indeed varies widely across different objects—for example, the "orange" and "lemon" 3D models depend much more on background than the "tennis ball." We also find that certain backgrounds yield systemically higher or lower accuracy; for example, average accuracy on "gray pier" is five times lower than that of "factory yard."

**Analyzing all backgrounds with the "coffee mug" model.** The previous study broadly characterizes classifier sensitivity classifiers to different models and environments. Now, to gain a deeper understanding of this sensitivity, we focus our analysis only a single 3D model (a "coffee mug") rendered in all 408 environments. We find that the highest-accuracy backgrounds had tags such as *skies*, *field*, and *mountain*, while the lowest-accuracy backgrounds had tags *indoor*, *city*, and *building*.

At first, this observation seems to be at odds with the idea that the classifier relies heavily on context clues to make decisions. After all, the backgrounds where the classifier seems to perform well (poorly) are places that we would expect a coffee mug to be rarely (frequently) present in the real world. Visualizing the best and worst backgrounds in terms of accuracy (Figure 6) suggests a possible explanation for this: the best backgrounds tend to be clean and distraction-free. Conversely, complicated backgrounds (e.g., some indoor scenes) often contain context clues that

---

[5]For computational reasons, we subsampled 20 environments which we used to analyze all of the pre-collected 3D models.

Figure 5: Examples of rendered scenes of the coffee mug 3D model in different environments, labeled with a pre-trained model's top prediction.
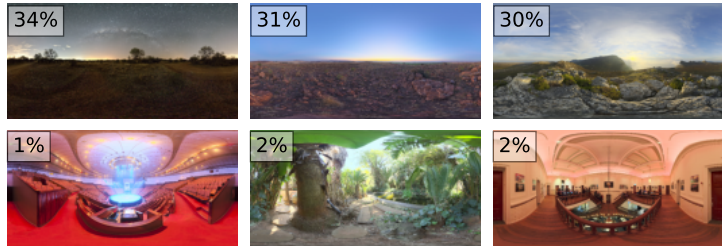


Figure 6: **(Top)** Best and **(Bottom)** worst background environments for classification of the coffee mug, and their respective accuracies (averaged over camera positions and zoom factors).
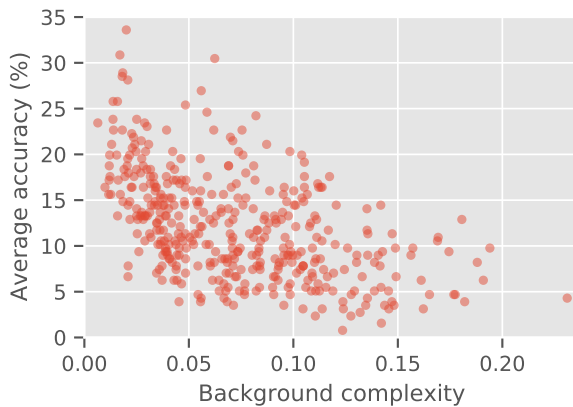


Figure 7: Relation between the complexity of a background and its average accuracy. Here complexity is defined as the average pixel value of the image after applying an edge detection filter.
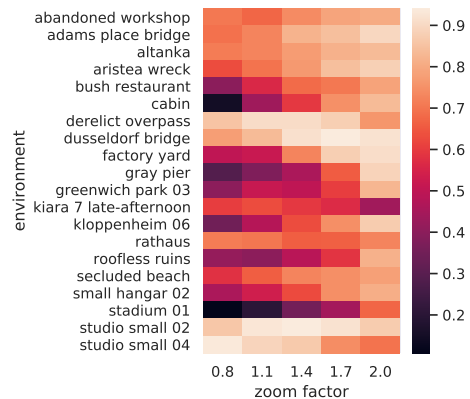


Figure 8: *3DB*'s focus on composability enables us to study robustness along multiple axes simultaneously. Here we study average model accuracy (computed over pose randomization) as a function of *both* zoom level and background.

make the mug difficult for models to detect. Comparing a "background complexity" metric (based on the number of edges in the image) to accuracy (Figure 7) supports this explanation: mugs overlaid on more complex backgrounds are more frequently misclassified by the model. In fact, some specific backgrounds even result in the model "hallucinating" objects; for example, the second-most frequent predictions for the *pond* and *sidewalk* backgrounds were *birdhouse* and *traffic light* respectively, despite the fact that neither object is present in the environment.

**Zoom/background interactions case study: the advantage of composable controls.** Finally, we leverage *3DB*'s composability to study interactions between controls. In Figure 8 we plot the mean classification accuracy of our "orange" model while varying background and scale factor. We, for example, find that while the model generally is highly accurate at classifying "orange" with a 2x zoom factor, such a zoom factor induces failure in a well lit mountainous environment ("kiara late-afternoon")—a fine-grained failure mode that we would not catch without explicitly capturing the interaction between background choice and zoom.

## 3.2  Texture-shape bias

We now demonstrate how *3DB* can be straightforwardly extended to discover more complex failure modes in computer vision models. Specifically, we will show how to rediscover the "texture bias" exhibited by ImageNet-trained

neural networks [GRM+19] in a systematic and (near-)photorealistic way. Geirhos et al. [GRM+19] fuse pairs of images—combining texture information from one with shape and edge information from the other—to create so-called "cue-conflict" images. They then demonstrate that on these images (cf. Figure 9), ImageNet-trained convolutional neural networks (CNNs) typically predict the class corresponding to the texture component, while humans typically predict based on shape features.
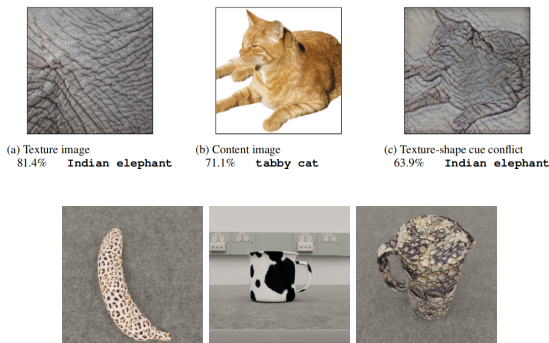


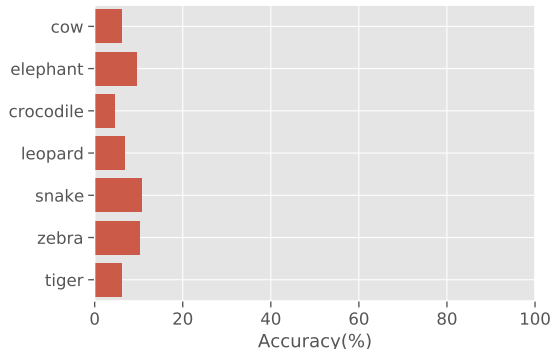Figure 9: Cue-conflict images generated by Geirhos et al. [GRM+19] (*top*) and *3DB* (*bottom*).



Figure 10: Model accuracy on previously correctly-classified images after their texture is altered via *3DB*, as a function of texture-type.

Cue-conflict images identify a concrete difference between human and CNN decision mechanisms. However, the fused images are unrealistic and can be cumbersome to generate (e.g., even the simplest approach uses style transfer [GEB16]). *3DB* gives us an opportunity to rediscover the influence of texture in a more streamlined fashion.

Specifically, we implement a control (now pre-packaged with *3DB*) that replaces an object's texture with a random (or user-specified) one. We use this control to create cue-conflict objects out of eight 3D models[6] and seven animal-skin texture images[7] (i.e., 56 objects in total). We test our pre-trained ResNet-18 on images of these objects rendered in a variety of poses and camera locations. Figure 9 displays sample cue-conflict images generated using *3DB*.

Our study confirms the findings of Geirhos et al. [GRM+19] and indicates that texture bias indeed extends to (near-)realistic settings. For images that were originally correctly classified (i.e., when rendered with the original texture), changing the texture reduced accuracy by 90-95% uniformly across textures (Figure 10). Furthermore, we observe that the model predictions usually align better with the texture of the objects rather than their geometry (Figure 11). One notable exception is the pitcher object, for which the most common prediction (aggregated over all textures) was *vase*. A possible explanation for this (based on inspection of the training data) is that due to high variability of vase textures in the train set, the classifier was forced to rely more on shape.

## 3.3 Orientation and scale dependence

Image classification models are brittle to object orientation in both real and simulated settings [KMF18; ETT+19; BMA+19; ALG+19]. As was the case for both background and texture sensitivity, reproducing and extending such observations is straightforward with *3DB*. Once again, we use the built-in controls to render objects at varying poses, orientations, scales, and environments before stratifying on properties of interest. Indeed, we find that classification accuracy is highly dependent on object orientation (Figure 13 left) and scale (Figure 13 right). However, this dependence is not uniform across objects. As one would expect, the classifier's accuracy is less sensitive to orientation on more symmetric objects (like "tennis ball" or "baseball"), but can vary widely on more uneven objects (like "drill").

For a more fine-grained look at the importance of object orientation, we can measure the classifier accuracy conditioned on a given part of each 3D model being visible. This analysis is once again straightforward in *3DB*, since each rendering is (optionally) accompanied by a UV map which maps pixels in the scene back to locations on on the object surface. Combining these UV maps with accuracy data allows one to construct the "accuracy heatmaps" shown in Figure 12, wherein each part of an object's surface corresponds to classifier accuracy on renderings in which the part is visible. The results confirm that atypical viewpoints adversely impact model performance, and

---

[6]Object models: mug, helmet, hammer, strawberry, teapot, pitcher, bowl, lemon, banana and spatula
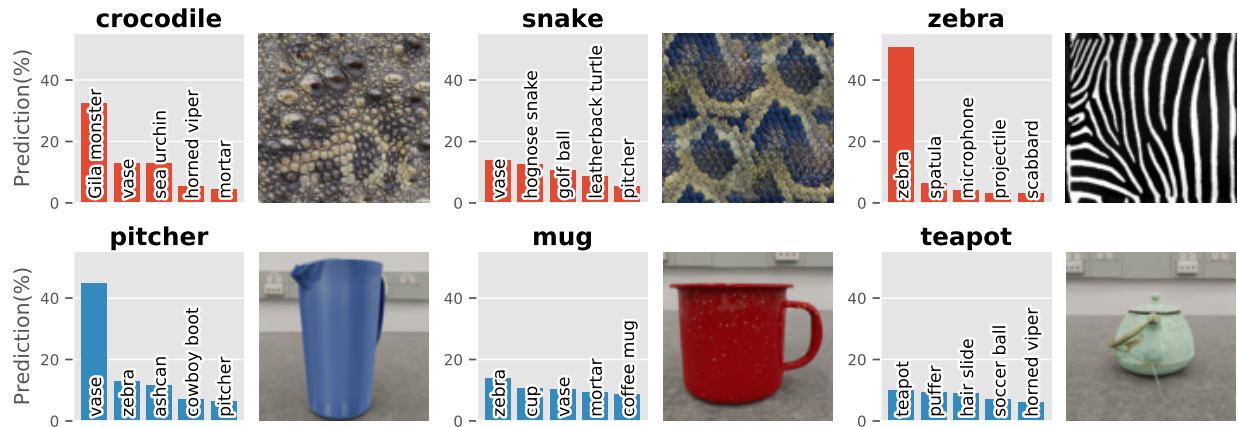[7]Texture types: cow, crocodile, elephant, leopard, snake, tiger and zebra

Figure 11: Distribution of classifier predictions after the texture of the 3D object model is altered. In the top row, we visualize the most frequently predicted classes for each texture (averaged over all objects). In the bottom row, we visualize the most frequently predicted classes for each object (averaged over all textures). We find that the model tends to predict based on the texture more often than based on the object.
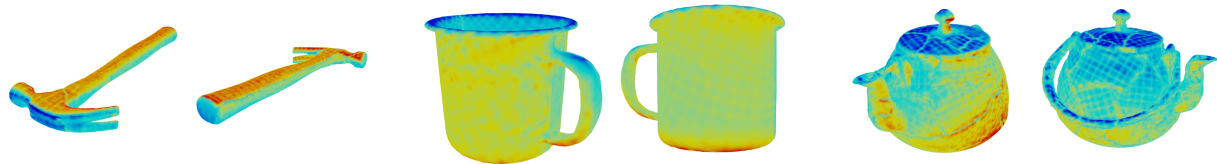


Figure 12: Model sensitivity to pose. The heatmaps denote the accuracy of the model in predicting the correct label, conditioned on a specific part of the object being visible in the image. Here, red and blue denotes high and low accuracy respectively.

also allow users to draw up a variety of testable hypotheses regarding performance on specific 3D models (e.g., for the coffee mug, the bottom rim is highlighted in red—is it the case that mugs are more accurately classified when viewed from the bottom)? These hypotheses can then be investigated further through natural data collection, or—as we discuss in the upcoming section—through additional experimentation with *3DB*.

## 3.4 Case study: using *3DB* to dive deeper

Our heatmap analysis in the previous section (cf. Figure 12) showed that classification accuracy for the mug decreases when its interior is visible. What could be causing this effect? One hypothesis is that in the ImageNet training set, objects are captured in context, and thus ImageNet-trained classifiers rely on this context to make decisions. Inspecting the ImageNet dataset, we notice that coffee mugs in context usually contain coffee in them. Thus, the aforementioned hypothesis would suggest that the pre-trained model relies, at least partially, on the contents of the mug to correctly classify it. *Can we leverage 3DB to confirm or refute this hypothesis?*

To test this, we implement a custom control that can render a liquid inside the "coffee mug" model. Specifically, this control takes water:milk:coffee ratios as parameters, then uses a parametric Blender shader (cf. Appendix F) to render a corresponding mixture of the liquids into the mug. We used the pre-packaged grid search policy, (programmatically) restricting the search space to viewpoints from which the interior of the mug was visible.

The results of the experiment are shown in Figure 14. It turns out that the model is indeed sensitive to changes in liquid, supporting our hypothesis: model predictions stayed constant (over all liquids) for only 20.7% of the rendered viewpoints (cf. Figure 14b). The *3DB* experiment provides further support for the hypothesis when we look at the correlation between the liquid mixture and the predicted class: Figure 14a visualizes this correlation in a normalized heatmap (for the unnormalized version, see Figure 22b in the Appendix F). We find that the model is most likely to predict "coffee mug" when coffee is added to the interior (unsurprisingly); as the coffee is mixed with water or
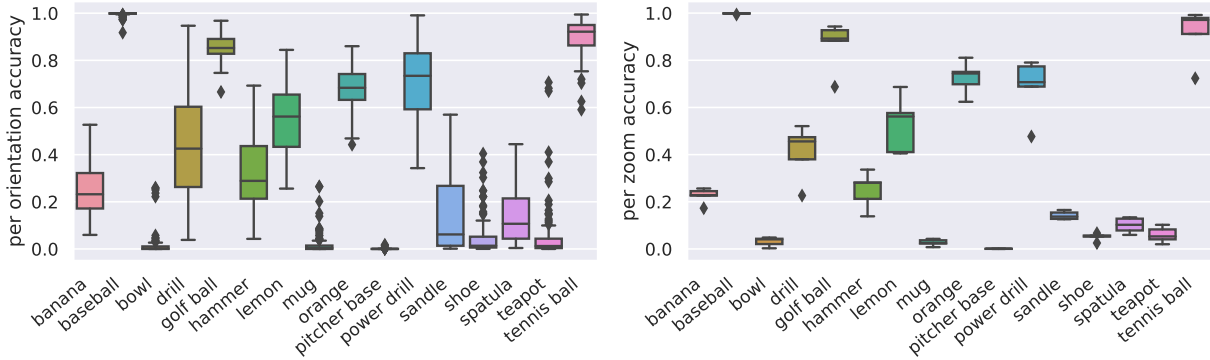
Figure 13: **(Left)** We compute the accuracy of the model for each object-orientation pair. For each object on the x-axis, we plot the variation in accuracy (over the set of possible orientations) using a boxplot. We visualize the per-orientation accuracy spread by including the median line, the first and third quartiles box edges, the range, and the outliers. **(Right)** Using the same format as the left hand plot, we plot how the classified object (on the x-axis) impacts variation in accuracy (over different zoom values) on the y-axis.
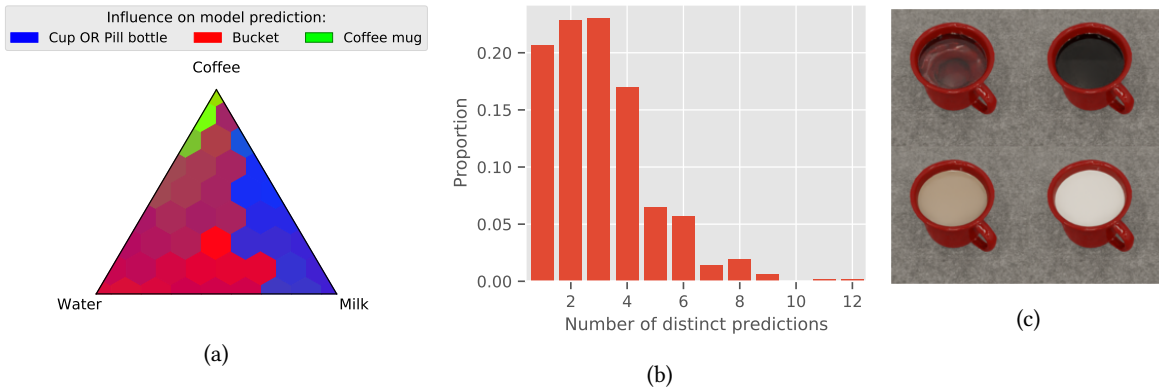


Figure 14: Testing classifier sensitivity to context: Figure (a) shows the correlation of the liquid mixture in the mug on the prediction of the model, averaged over random viewpoints (see Figure 22b for the raw frequencies). Figure (b) shows that for a fixed viewpoint, model predictions are unstable with respect to the liquid mixture. Figure (c) shows examples of rendered liquids (water, black coffee, milk, and milk/coffee mix).

milk, the predicted label distribution shifts towards "bucket" and "cup" or "pill bottle," respectively. Overall, our experiment suggests that current ResNet-18 classifiers are indeed sensitive to object context—in this case, the fluid composition of the mug interior. More broadly, this illustration highlights how a system designer can quickly go from hypothesis to empirical verification with minimal effort using *3DB*. (In fact, going from the initial hypothesis to Figure 14 took less than a single day of work for one author.)

# 4   Physical realism

The previous sections have demonstrated various ways in which we can use *3DB* to obtain insights into model behavior in simulation. Our overarching goal, however, is to understand when models will fail in the physical world. Thus, we would like for the insights extracted by *3DB* to correspond to naturally-arising model behavior, and not just artifacts of the simulation itself [8]. To this end, we now test the *physical realism* of *3DB*: can we understand model

---

[8]Indeed, a related challenge is the *sim2real* problem in reinforcement learning, where agents trained in simulation latch on to simulator properties and fail to generalize to the real world. In both cases, we are concerned about artifacts or spurious correlations that invalidate conclusions made in simulation.
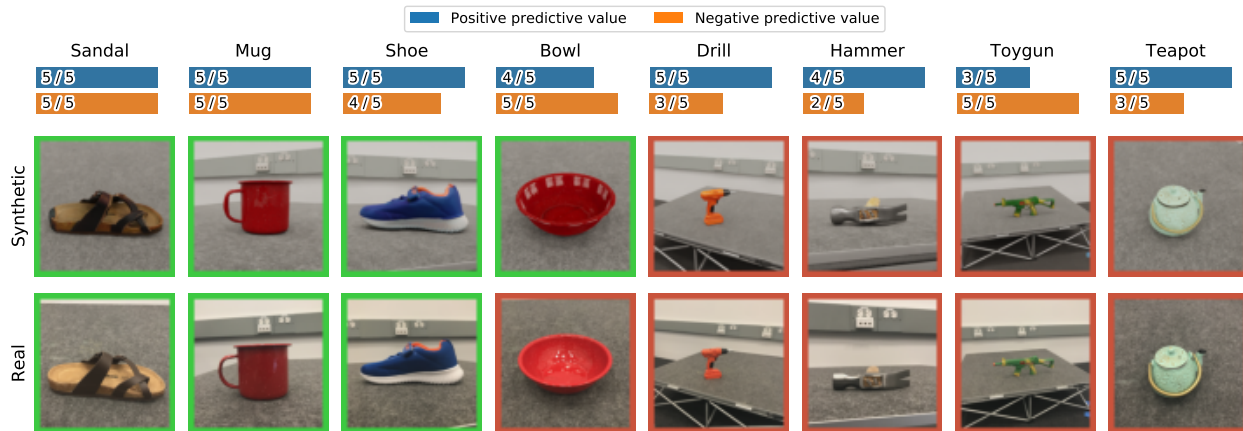
Figure 15: **(Top)** Agreement, in terms of model correctness, between model predictions within *3DB* and model predictions in the real world. For each object, we selected five rendered scenes found by *3DB* that were misclassified in simulation, and five that were correctly classified; we recreated and deployed the model on each scene in the physical world. The *positive (resp., negative) predictive value* is rate at which correctly (resp. incorrectly) classified examples in simulation were also correctly (resp., incorrectly) classified in the physical world. **(Bottom)** Comparison between example simulated scenes generated by *3DB* (first row) and their recreated physical counterparts (second row). Border color indicates whether the model was correct on this specific image.

performance (and uncover vulnerabilities) on real photos using only a high-fidelity simulation?

To answer this question, we collected a set of physical objects with corresponding 3D models, and set up a physical room with its corresponding 3D environment. We used *3DB* to identify strong points and vulnerabilities of a pre-trained ImageNet classifier in this environment, mirroring our methodology from Section 3. We then recreated each scenario found by *3DB* in the physical room, and took photographs that matched the simulation as closely as possible. Finally, we evaluated the physical realism of the system by comparing models' performance on the photos (i.e., whether they classified each photo correctly) to what *3DB* predicted.

**Setup.** We performed the experiment in the studio room shown in Appendix Figure 20b for which we obtained a fairly accurate 3D model (cf. Appendix Figure 20a). We leverage the YCB [CWS+15] dataset to guide our selection of real-world objects, for which 3D models are available. We supplement these by sourcing additional objects (from amazon.com) and using a 3D scanner to obtain corresponding meshes. [9]

We next used *3DB* to analyze the performance of a pre-trained ImageNet ResNet-18 on the collected objects in simulation, varying over a set of realistic object poses, locations, and orientations. For each object, we selected 10 rendered situations: five where the model made the correct prediction, and five where the model predicted incorrectly. We then tried to recreate each rendering in the physical world. First we roughly placed the main object in the location and orientation specified in the rendering, then we used a custom-built iOS application (see Appendix C) to more precisely match the rendering with the physical setup.

**Results.** Figure 15 visualizes a few samples of renderings with their recreated physical counterparts, annotated with model correctness. Overall, we found a 85% agreement rate between the model's correctness on the real photos and the synthetic renderings—agreement rates per class are shown in Figure 15. Thus, despite imperfections in our physical reconstructions, the vulnerabilities identified by *3DB* turned out to be physically realizable vulnerabilities (and conversely, the positive examples found by *3DB* are usually also classified correctly in the real world). We found that objects with simpler/non-metallic materials (e.g., the bowl, mug, and sandal) tended to be more reliable than metallic objects such as the hammer and drill. It is thus possible that more precise texture tuning of 3D models object could increase agreement further (although a more comprehensive study would be needed to verify this).

---

[9]We manually adjusted the textures of these 3D models to increase realism (e.g., by tuning reflectance or roughness). In particular, classic photogrammetry is unable to model the metallicness and reflectivity of objects. It also tends to embed reflections as part of the color of the object

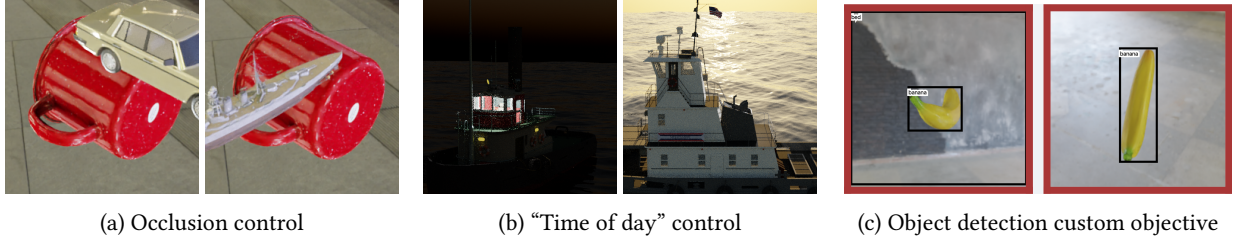| (a) Occlusion control | (b) "Time of day" control | (c) Object detection custom objective |

Figure 16: Example of some of the ways in which one can extend *3DB*: adding custom controls, defining custom objectives, and integrating external libraries.

# 5 Extensibility

*3DB* was designed with extensibility in mind. Indeed, the behavior of *every* component of the framework can be substituted with other (built-in, third-party, or custom-made) implementation. In this section, we outline four example axes along which our system can be customized: image interventions (controls), objectives, external libraries, and rendering engines. Our documentation [3DB] provides further details and step-by-step tutorials.

**Custom controls.** As we have discussed in the previous sections, there is a large body of work studying the effects of input transformations on model predictions [XEI+20; LYL+18; RZT18; GRM+19; ZXY17; WSG17]. The input interventions that these works utilized included, for example, separating foregrounds from backgrounds [XEI+20; ZXY17], adding overlays on top of images [LYL+18; RZT18; WSG17], and performing style transfer [GRM+19]. These interventions have been implemented with a lot of care. However, they still tend to introduce artifacts and can lack realism. In Section 3 we already demonstrated that *3DB* is able to circumvent these problems in a streamlined and composable manner. Indeed, by operating in three dimensional space, i.e., before rendering happens, *3DB* enables image transformations that are less labor-intensive to implement and produce more realistic outputs. To showcase this, in Section 3 we replicated various image transformation studies using the *controls* built in to *3DB* (e.g., Figure 10 corresponds to the study of [GRM+19]). However, beyond these built-in capabilities, users can also add custom controls that implement their desired transformations: Figure 16a, for example, depicts the output of a custom "occlusion control" that could be used to replicate studies such as [RZT18].

**Custom objectives.** Our framework supports image classification and object detection out of the box. (In this paper, we focus primarily on the former—cf. Figure 16c for an example of the latter.) Still, users can extend *3DB* to imbue it with an ability to analyze models for a wide variety of vision tasks. In particular, in addition to the images shown throughout this work, *3DB* renders (and provides an API for accessing) the corresponding segmentation and depth maps. This allow users to easily use the framework for tasks such as depth estimation, instance segmentation, and image segmentation (the last one of these is in fact subject of our tutorial on the implementation of custom tasks[10]). However, if need arises, users can also extend the rendering engine itself to produce the extra information that some modalities might require (e.g., the coordinates of joints for pose estimation).

**External libraries.** *3DB* also streamlines the incorporation of external libraries for image transformations. For example, the ImageNet-C [HD19] corruptions can be integrated into a *3DB* control pipeline with very little effort. (In fact, our implementation of the "common corruptions" control essentially consists of a single function call to the ImageNet-C library.)

**Rendering engine.** Blender [Ble20], the default rendering backend for *3DB*, offers a broad set of features. Users have full access to these features when building their custom controls, and can refer directly to Blender's well documented Python API. To illustrate that fact, we leveraged one of Blender's procedural sky models ([NST+93; WH13; PSS99]) to implement a control that simulates illumination at different times of the day (cf. Figure 16b).

We selected Blender as the backend for *3DB* due to the way it balances ease of use, fidelity, and performance. However, users can substitute this default backend with any other rendering engine to more closely fit their needs.

---

[10]https://3db.github.io/3db/usage/custom_evaluator.html

For example, users can, on the one hand, setup a rendering backend (and corresponding controls) based on Mitsuba [NVZ+19], a research-oriented engine capable of highly accurate simulation. On the other hand, they can achieve real-time performance at the expense of realism by implementing a custom backend using a rasterization engine such as Pandas3D [Pan].

# 6   Related Work

*3DB* builds on a growing body of work that looks beyond accuracy-based benchmarks in order to understand the *robustness* of modern computer vision models and their failure modes. In particular, our goal is to provide a unified framework for reproducing these studies and for conducting new such analyses. In this section, we discuss the existing research in robustness, interpretability, and simulation that provide the context for our work.

**Adversarial robustness.**   Several recent works propose analyzing model robustness by crafting adversarial, i.e., *worst-case*, inputs. For example, [SZS+14] discovered that a carefully chosen but imperceptible perturbation suffices to change classifier predictions on virtually any natural input. Subsequently, the study of such "adversarial examples" has extended far beyond the domain of image classification: e.g., recent works have studied worst-case inputs for object detection and image segmentation [EEF+18; XWZ+17; FKM+17]; generative models [KFS18]; and reinforcement learning [HPG+17]. More closely related to our work are studies focused on three-dimensional or physical-world adversarial examples [EEF+18; BMR+18; AEI+18; XYL+19; LTL+19]. These studies typically use differentiable rendering and perturb object texture, geometry, or lighting to induce misclassification. Alternatively, Li, Schmidt, and Kolter [LSK19] modify the camera itself via an adversarial camera lens that consistently cause models to misclassify inputs.

In our work, we have primarily focused on using non-differentiable but high-fidelity rendering to analyze a more *average-case* notion of model robustness to semantic properties such as object orientation or image backgrounds. Nevertheless, the extensibility of *3DB* means that users can reproduce such studies (by swapping out the Blender rendering module for a differentiable renderer, writing a custom control, and designing a custom search policy) and use our framework to attain a more realistic understanding of the worst-case robustness of vision models.

**Robustness to synthetic perturbations.**   Another popular approach to analyzing model robustness involves applying transformations to natural images and measuring the resultant changes in model predictions. For example, Engstrom et al. [ETT+19] measures robustness to image rotations and translations; Geirhos et al. [GRM+19] study robustness to style transfer (i.e., texture perturbations); and a number of works has studied robustness to common corruptions [HD19; KSH+19], changes in image backgrounds [ZXY17; XEI+20], Gaussian noise [FGC+19], and object occlusions [RZT18], among other transformations.

A more closely related approach to ours analyzes the impact of factors such as object pose and geometry by applying synthetic perturbations in three-dimensional space [HG19; SLQ+20; HMG18; ALG+19]. For example, Hamdi and Ghanem [HG19] and Jain et al. [JCJ+20] use a neural mesh renderer [KUH18] and Redner [LAD+18] respectively to render images to analyze the failure modes of vision models. Alcorn et al. [ALG+19] present a system for discovering neural networks failure modes as a function of object orientation, zoom, and (two-dimensional) background and perform a thorough study on the impact of these factors on model decisions.

*3DB* draws inspiration from the studies listed above and tries to provide a unified framework for detecting *arbitrary* model failure modes. For example, our framework provides explicit mechanisms for users to make custom controls and custom search strategies, and includes built-in controls designed to range many possible failure modes encompassing nearly all of the aforementioned studies (cf. Section 3). Users can also *compose* different transformations in *3DB* to get an even more fine-grained understanding of model robustness.

**Other types of robustness.**   An oft-studied but less related branch of robustness research tests model performance on unaltered images from distributions that are nearby but non-identical to that of the training set. Examples of such investigations include studies of newly collected datasets such as ImageNet-v2 [RRS+19; EIS+20; TDS+20], ObjectNet [BMA+19], and others (e.g., [HZB+19; SDR+19]). In a similar vein, Torralba and Efros [TE11] study model performance when trained on one standard dataset and tested on another. We omit a detailed discussion of these works since *3DB* is synthetic by nature (and thus less photorealistic than the aforementioned studies). As shown in Section 4, however, *3DB* is indeed realistic enough to be indicative of real-world performance.

**Interpretability, counterfactuals, and model debugging.** *3DB* can be cast as a method for *debugging* vision models that provides users fine-grained control over the rendered scenes and thus enabling them to find specific modes of failure (cf. Sections 3 and 4). Model debugging is also a common goal in intepretability research, where methods generally seek to provide justification for model decisions based on either local features (i.e., specific to the image at hand) or global ones (i.e., general biases of the model). Local explanation methods, including saliency maps [SVZ13; DG17; STY17], surrogate models such as LIME [RSG16], and counterfactual image pairs [FV17; ZXY17; GWE+19], can provide insight into specific model decisions but can also be fragile [GAZ19; AJ18] or misleading with respect to global model behaviour [AGM+18; STY17; AML+20; Lip18]. Global interpretability methods include concept-based explanations [BZK+17; KWG+18; YKA+20; WSM21] (though such explanations can often lack causal link to the features models actually use [GWE+19]), but also encompass many of the robustness studies highlighted earlier in this section, which can be cast as uncovering global biases of vision models.

**Simulated environments and training data.** Finally, there has been a long line of work on developing simulation platforms that can serve as both a source of additional (synthetic) training data, and as a proxy for real-world experimentation. Such simulation environments are thus increasingly playing a role in fields such as computer vision, robotics, and reinforcement learning (RL). For instance, OpenAI Gym [BCP+16] and DeepMind Lab [BLT+16] provide simulated RL training environments with a fleet of control tasks. Other frameworks such as UnityML [JBT+20] and RoboSuite [ZWM+20] were subsequently developed to cater to more complex agent behavior.

In computer vision, the Blender rendering engine [Ble20] has been used to generate synthetic training data through projects such as BlenderProc [DSW+19] and BlendTorch [HBZ+20]. Similarly, HyperSim [RP] is a photorealistic synthetic dataset focused on multimodal scene understanding. Another line of work learns optimal simulation parameters for synthetic data generation according to user-defined objective, such as minimizing the distribution gap between train and test environments [KPL+19; DKF20; BBG+20]. Simulators such as AirSim [SDL+18], FlightMare [SNK+20], and CARLA [DRC+17] (built on top of video game engines Unreal Engine and Unity) allow for collection of synthetic training data for perception and control. In robotics, simulators include environments that model typical household layouts for robot navigation [KMH+17; WWG+18; PRB+18], interactive ones where objects that can be actuated [XZH+18; XSL+20; XQM+20], and those that include support for tasks such as question answering and instruction following [SKM+19].

While some of these platforms may share components with *3DB* (e.g., the physics engine, photorealistic rendering), the do not share the same goals as *3DB*, i.e., diagnosing specific failures in existing models.

# 7 Conclusion

In this work, we introduced *3DB*, a unified framework for diagnosing failure modes in vision models based on high-fidelity rendering. We demonstrate the utility of *3DB* by applying it to a number of model debugging use cases—such as understanding classifier sensitivities to realistic scene and object perturbations, and discovering model biases. Further, we show that the debugging analysis done using *3DB* in simulation is actually predictive of model behavior in the physical world. Finally, we note that *3DB* was designed with extensibility as a priority; we encourage the community to build upon the framework so as to uncover new insights into the vulnerabilities of vision models.

# Acknowledgements

# References

[3DB]      3DB. *Documentation*. URL: https://3db.github.io/3db/.

[AEI+18]   Anish Athalye et al. "Synthesizing Robust Adversarial Examples". In: *International Conference on Machine Learning (ICML)*. 2018.

[AGM+18]   Julius Adebayo et al. "Sanity checks for saliency maps". In: *Neural Information Processing Systems (NeurIPS)*. 2018.

[AJ18]     David Alvarez-Melis and Tommi S Jaakkola. "On the robustness of interpretability methods". In: *arXiv preprint arXiv:1806.08049* (2018).

[ALG+19]   Michael A Alcorn et al. "Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects". In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019.

[AML+20]   Julius Adebayo et al. "Debugging Tests for Model Explanations". In: 2020.

[BBG+20]   Harkirat Singh Behl et al. "Autosimulate:(quickly) learning synthetic data generation". In: *European Conference on Computer Vision*. Springer. 2020, pp. 255–271.

[BCP+16]   Greg Brockman et al. "Openai gym". In: *arXiv preprint arXiv:1606.01540* (2016).

[Ble20]    Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation. Stichting Blender Foundation, Amsterdam, 2020. URL: http://www.blender.org.

[BLT+16]   Charles Beattie et al. "Deepmind lab". In: *arXiv preprint arXiv:1612.03801* (2016).

[BMA+19]   Andrei Barbu et al. "ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models". In: *Neural Information Processing Systems (NeurIPS)*. 2019.

[BMR+18]   Tom B. Brown et al. *Adversarial Patch*. 2018. arXiv: 1712.09665 [cs.CV].

[BZK+17]   David Bau et al. "Network dissection: Quantifying interpretability of deep visual representations". In: *Computer Vision and Pattern Recognition (CVPR)*. 2017.

[CWS+15]   Berk Calli et al. "Benchmarking in manipulation research: The YCB object and model set and benchmarking protocols". In: *arXiv preprint arXiv:1502.03143* (2015).

[DG17]     Piotr Dabkowski and Yarin Gal. "Real time image saliency for black box classifiers". In: *Neural Information Processing Systems (NeurIPS)*. 2017.

[DKF20]    Jeevan Devaranjan, Amlan Kar, and Sanja Fidler. "Meta-Sim2: Unsupervised Learning of Scene Structure for Synthetic Data Generation". In: *European Conference on Computer Vision*. Springer. 2020, pp. 715–733.

[DRC+17]   Alexey Dosovitskiy et al. "CARLA: An open urban driving simulator". In: *arXiv preprint arXiv:1711.03938* (2017).

[DSW+19]   Maximilian Denninger et al. "BlenderProc". In: *arXiv preprint arXiv:1911.01911* (2019).

[EEF+18]   Kevin Eykholt et al. "Physical Adversarial Examples for Object Detectors". In: *CoRR* (2018).

[EIS+20]   Logan Engstrom et al. "Identifying Statistical Bias in Dataset Replication". In: *International Conference on Machine Learning (ICML)*. 2020.

[ETT+19]   Logan Engstrom et al. "Exploring the Landscape of Spatial Robustness". In: *International Conference on Machine Learning (ICML)*. 2019.

[FGC+19]   Nic Ford et al. "Adversarial Examples Are a Natural Consequence of Test Error in Noise". In: *arXiv preprint arXiv:1901.10513*. 2019.

[FKM+17]   Volker Fischer et al. "Adversarial examples for semantic image segmentation". In: *Arxiv preprint arXiv:1703.01101*. 2017.

[FV17]     Ruth C Fong and Andrea Vedaldi. "Interpretable explanations of black boxes by meaningful perturbation". In: *International Conference on Computer Vision (ICCV)*. 2017.

[GAZ19]    Amirata Ghorbani, Abubakar Abid, and James Zou. "Interpretation of neural networks is fragile". In: *AAAI Conference on Artificial Intelligence (AAAI)*. 2019.
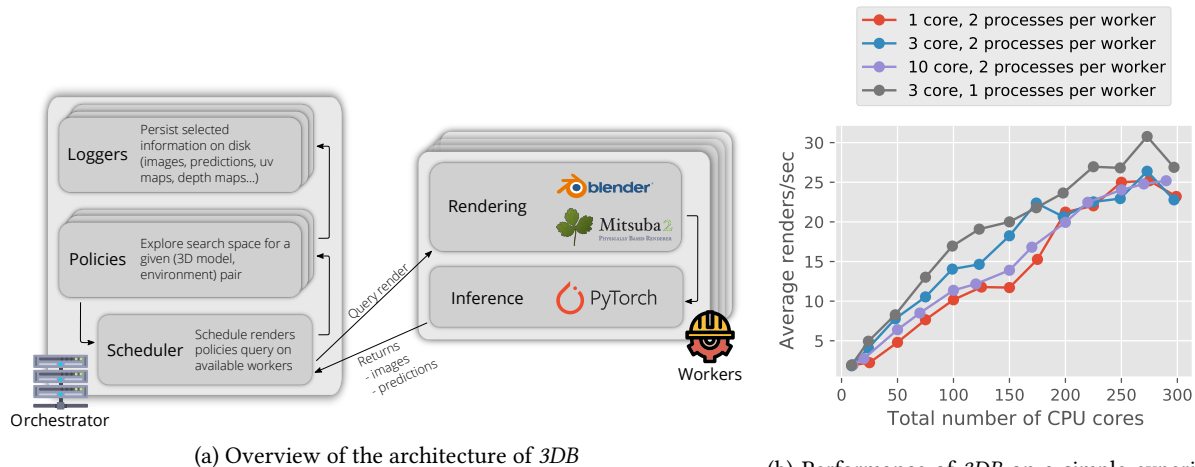
[GEB16]    Leon A Gatys, Alexander S Ecker, and Matthias Bethge. "Image style transfer using convolutional neural networks". In: *computer vision and pattern recognition (CVPR)*. 2016.

[GRM+19]   Robert Geirhos et al. "ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness." In: *International Conference on Learning Representations (ICLR)*. 2019.

[GWE+19]   Yash Goyal et al. "Counterfactual visual explanations". In: *arXiv preprint arXiv:1904.07451* (2019).

[HBZ+20]   Christoph Heindl et al. "BlendTorch: A Real-Time, Adaptive Domain Randomization Library". In: *arXiv preprint arXiv:2010.11696* (2020).

[HD19]     Dan Hendrycks and Thomas G. Dietterich. "Benchmarking Neural Network Robustness to Common Corruptions and Surface Variations". In: *International Conference on Learning Representations (ICLR)*. 2019.

[HG19]     Abdullah Hamdi and Bernard Ghanem. "Towards Analyzing Semantic Robustness of Deep Neural Networks". In: *arXiv preprint arXiv:1904.04621* (2019).

[HMG18]    Abdullah Hamdi, Matthias Muller, and Bernard Ghanem. "SADA: Semantic Adversarial Diagnostic Attacks for Autonomous Applications". In: *arXiv preprint arXiv:1812.02132* (2018).

[HPG+17]   Sandy Huang et al. "Adversarial Attacks on Neural Network Policies". In: *ArXiv preprint arXiv:1702.02284*. 2017.

[HZB+19]   Dan Hendrycks et al. "Natural adversarial examples". In: *arXiv preprint arXiv:1907.07174* (2019).

[HZR+15]   Kaiming He et al. *Deep Residual Learning for Image Recognition*. 2015.

[JBT+20]   Arthur Juliani et al. *Unity: A General Platform for Intelligent Agents*. 2020. arXiv: 1809.02627 [cs.LG].

[JCJ+20]   Lakshya Jain et al. "Analyzing and Improving Neural Networks by Generating Semantic Counterexamples through Differentiable Rendering". In: *arXiv preprint arXiv:1910.00727* (2020).

[KFS18]    Jernej Kos, Ian Fischer, and Dawn Song. "Adversarial examples for generative models". In: *IEEE Security and Privacy Workshops (SPW)*. 2018.

[KMF18]    Can Kanbak, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. "Geometric robustness of deep networks: analysis and improvement". In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2018.

[KMH+17]   Eric Kolve et al. "Ai2-thor: An interactive 3d environment for visual ai". In: *arXiv preprint arXiv:1712.05474* (2017).

[KPL+19]   Amlan Kar et al. "Meta-sim: Learning to generate synthetic datasets". In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019, pp. 4551–4560.

[KSH+19]   Daniel Kang et al. "Testing Robustness Against Unforeseen Adversaries". In: *ArXiv preprint arxiv:1908.08016*. 2019.

[KUH18]    Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. "Neural 3D Mesh Renderer". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2018.

[KWG+18]   Been Kim et al. "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav)". In: *International conference on machine learning (ICML)*. 2018.

[LAD+18]   Tzu-Mao Li et al. "Differentiable Monte Carlo Ray Tracing through Edge Sampling". In: *SIGGRAPH Asia 2018 Technical Papers*. 2018.

[Lip18]    Zachary C Lipton. "The Mythos of Model Interpretability: In machine learning, the concept of interpretability is both important and slippery." In: (2018).

[LSK19]    Juncheng Li, Frank R. Schmidt, and J. Zico Kolter. "Adversarial camera stickers: A physical camera-based attack on deep learning systems". In: *Arxiv preprint arXiv:1904.00759*. 2019.

[LTL+19]   Hsueh-Ti Derek Liu et al. "Beyond Pixel Norm-Balls: Parametric Adversaries Using An Analytically Differentiable Renderer". In: *International Conference on Learning Representations (ICLR)*. 2019.

[LYL+18]   Xin Liu et al. "Dpatch: An adversarial patch attack on object detectors". In: *arXiv preprint arXiv:1806.02299* (2018).

[NST+93] Tomoyuki Nishita et al. "Display of the Earth Taking into Account Atmospheric Scattering". In: *Proceedings of the 20th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. 1993.

[NVZ+19] Merlin Nimier-David et al. "Mitsuba 2: A Retargetable Forward and Inverse Renderer". In: *Transactions on Graphics (Proceedings of SIGGRAPH Asia)* 38.6 (2019). DOI: 10.1145/3355089.3356498.

[Pan] Pandas3D. *The Open Source Framework for 3D Rendering and Games*. URL: https://www.panda3d.org/.

[PRB+18] Xavier Puig et al. "Virtualhome: Simulating household activities via programs". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018.

[PSS99] A. J. Preetham, Peter Shirley, and Brian Smits. "A Practical Analytic Model for Daylight". In: *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. 1999.

[RDS+15] Olga Russakovsky et al. "ImageNet Large Scale Visual Recognition Challenge". In: *International Journal of Computer Vision (IJCV)*. 2015.

[RP] Mike Roberts and Nathan Paczan. *Hypersim: A Photorealistic Synthetic Dataset for Holistic Indoor Scene Understanding*. arXiv 2020.

[RRS+19] Benjamin Recht et al. "Do ImageNet Classifiers Generalize to ImageNet?" In: *International Conference on Machine Learning (ICML)*. 2019.

[RSG16] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ""Why should I trust you?" Explaining the predictions of any classifier". In: *International Conference on Knowledge Discovery and Data Mining (KDD)*. 2016.

[RZT18] Amir Rosenfeld, Richard Zemel, and John K. Tsotsos. "The Elephant in the Room". In: *arXiv preprint arXiv:1808.03305*. 2018.

[SDL+18] Shital Shah et al. "Airsim: High-fidelity visual and physical simulation for autonomous vehicles". In: *Field and service robotics*. Springer. 2018, pp. 621–635.

[SDR+19] Vaishaal Shankar et al. "Do Image Classifiers Generalize Across Time?" In: *arXiv preprint arXiv:1906.02168* (2019).

[SKM+19] Manolis Savva et al. "Habitat: A platform for embodied ai research". In: *Proceedings of the IEEE International Conference on Computer Vision*. 2019.

[SLQ+20] Michelle Shu et al. "Identifying Model Weakness with Adversarial Examiner". In: *AAAI Conference on Artificial Intelligence (AAAI)*. 2020.

[SNK+20] Yunlong Song et al. "Flightmare: A Flexible Quadrotor Simulator". In: *arXiv preprint arXiv:2009.00563* (2020).

[STY17] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. "Axiomatic attribution for deep networks". In: *International Conference on Machine Learning (ICML)*. 2017.

[SVZ13] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. "Deep inside convolutional networks: Visualising image classification models and saliency maps". In: *arXiv preprint arXiv:1312.6034* (2013).

[SZS+14] Christian Szegedy et al. "Intriguing properties of neural networks". In: *International Conference on Learning Representations (ICLR)*. 2014.

[TDS+20] Rohan Taori et al. "Measuring Robustness to Natural Distribution Shifts in Image Classification". In: *arXiv preprint arXiv:2007.00644* (2020).

[TE11] Antonio Torralba and Alexei A Efros. "Unbiased look at dataset bias". In: *CVPR 2011*. 2011.

[WH13] Alexander Wilkie and Lukas Hosek. "Predicting Sky Dome Appearance on Earth-like Extrasolar Worlds". In: *29th Spring conference on Computer Graphics (SCCG 2013)*. 2013.

[WSG17] Xiaolong Wang, Abhinav Shrivastava, and Abhinav Gupta. "A-Fast-RCNN: Hard Positive Generation via Adversary for Object Detection". In: *ArXiv preprint arXiv:1704.03414*. 2017.

[WSM21] Eric Wong, Shibani Santurkar, and Aleksander Madry. "Leveraging Sparse Linear Layers for Debuggable Deep Networks". In: *Arxiv preprint arXiv:2105.04857*. 2021.

[WWG+18]     Yi Wu et al. "Building generalizable agents with a realistic and rich 3d environment". In: *arXiv preprint arXiv:1801.02209* (2018).

[XEI+20]     Kai Xiao et al. "Noise or signal: The role of image backgrounds in object recognition". In: *arXiv preprint arXiv:2006.09994* (2020).

[XQM+20]     Fanbo Xiang et al. "SAPIEN: A simulated part-based interactive environment". In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 2020.

[XSL+20]     Fei Xia et al. "Interactive Gibson Benchmark: A Benchmark for Interactive Navigation in Cluttered Environments". In: *IEEE Robotics and Automation Letters* (2020).

[XWZ+17]     Cihang Xie et al. "Adversarial examples for semantic segmentation and object detection". In: *Proceedings of the IEEE International Conference on Computer Vision.* 2017, pp. 1369–1378.

[XYL+19]     Chaowei Xiao et al. "MeshAdv: Adversarial Meshes for Visual Recognition". In: *Computer Vision and Pattern Recognition (CVPR).* 2019.

[XZH+18]     Fei Xia et al. "Gibson env: Real-world perception for embodied agents". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* 2018.

[YKA+20]     Chih-Kuan Yeh et al. "On Completeness-aware Concept-Based Explanations in Deep Neural Networks". In: *Advances in Neural Information Processing Systems (NeurIPS)* (2020).

[ZML+07]     Jianguo Zhang et al. "Local features and kernels for classification of texture and object categories: A comprehensive study". In: *International journal of computer vision.* 2007.

[ZWM+20]     Yuke Zhu et al. "robosuite: A modular simulation framework and benchmark for robot learning". In: *arXiv preprint arXiv:2009.12293* (2020).

[ZXY17]      Zhuotun Zhu, Lingxi Xie, and Alan Yuille. "Object Recognition without and without Objects". In: *International Joint Conference on Artificial Intelligence.* 2017.

# A   Implementation and scalability

In this section, we briefly describe the underlying architecture of *3DB*, and verify that the system can effectively scale to distributed compute infrastructure.



(a) Overview of the architecture of *3DB*



(b) Performance of *3DB* on a simple experiment in function of the number of CPU cores recruited for renders.

**Architecture.**   To ensure scalability of this pipeline, we implement *3DB* as a client-server application (cf. Figure 17a). The main "orchestrator" thread constructs a search space by composing the user's specified controls, then uses the (user-specified) policy to find the exact set of 3D configurations that need to be rendered and analyzed. It then schedules these configurations across a set of worker nodes, whose job is to receive configurations, render them, run inference using the user's pretrained vision model, and send the results back to the orchestrator node. The results are aggregated and written to disk by a logging module. The dashboard is implemented as a separate entity that reads the log files and produces a user-friendly web interface for understanding the *3DB* results.

**Scalability.**   As discussed in Section 2, in order to perform photo-realistic rendering at scale, *3DB* must be able to leverage many machines (CPU cores) in parallel. *3DB* is designed to allow for this. It can accommodate as many rendering clients as the user can afford and the rendering efficiency of *3DB* largely scales linearly with available CPU cores (cf. Figure 17b). Note that the although the user can add as much rendering clients as they want, the number of actually used clients by the orchestrator is limited by its number of policy instances. In our paper, we run a limited number instances of policies (one instance per (env, 3D model) pair) concurrently to keep the memory of the orchestrator under control. This limits the scalability of the system as the maximum of renders that has to be done at any point in time scales with the number of policies of the orchestrator. Yet, were able to reach 415 FPS average/800 FPS peak throughput with dummy workers (no rendering), and around 100 FPS for the main experiments of this paper (e.g. physical realism experiment) which uses a complex background environment requiring substantial amount of rendering time (15 secs per image).
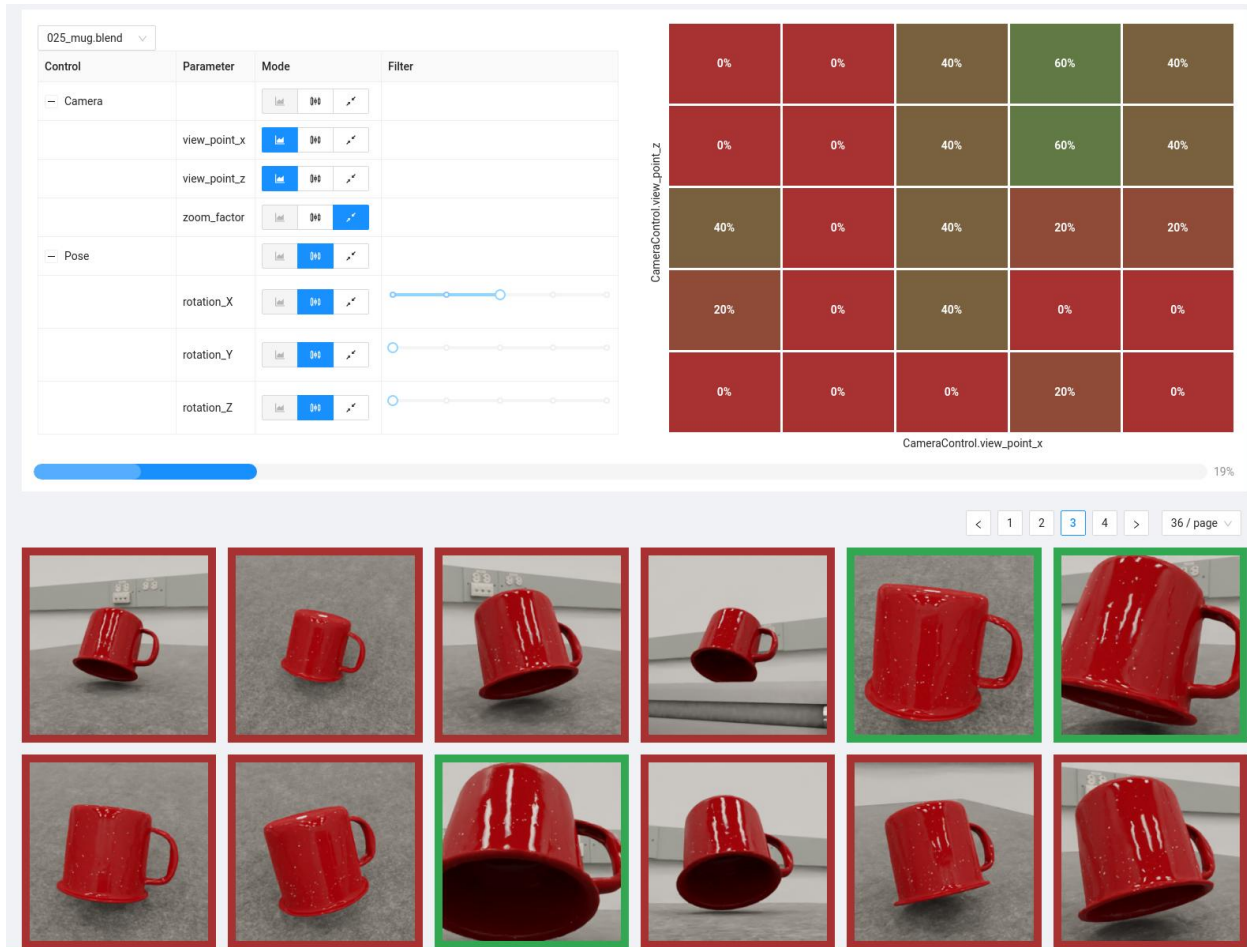
# B   Experiment Dashboard



Figure 18: Screenshot of the dashboard used for data exploration.

Since experiments usually produce large amounts of data that can be hard to get a sense of, we created a data visualization dashboard. Given a folder containing the JSON logs of a job, it offers a user interface to explore the influence of the controls.

For each parameter of each control, we can pick one out three mode:

- **Heat map axis**: This control will be used as the x or y axis of the heat map. Exactly two controls should be assigned to this mode to enable the visualization. Hovering on cells of the heat map will filter all samples falling in that region.

- **Slider**: This mode enables a slider that is used to only select the samples that match exactly this particular value.

- **Aggregate**: do not filter samples based on this parameter

# C   iPhone App

We developed a native iOS app to help align objects in the physical experiment (Section 4). The app allows the user to enter one or more rendering IDs (corresponding to scenes rendered by *3DB*); the app then brings up a camera with a translucent overlay of either the scene or an edge-filtered version of the scene (cf. Figure 19). We used the app to

align the physical object and environment with their intended place in the rendered scene. The app connects to the same backend serving the experiment dashboard.
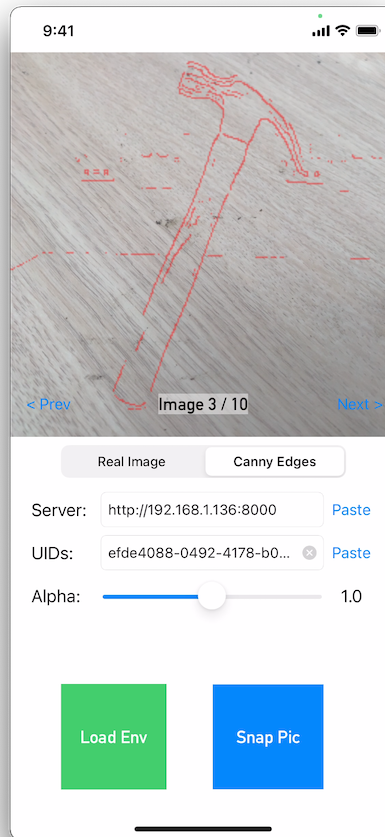


Figure 19: A screenshot of the iOS app used to align objects for the physical-world experiment. After starting the dashboard server, the user can specify the server location as well as a set of rendering IDs. The corresponding renderings will be displayed over a camera view, allowing the user to correctly position the object in the frame. The user can adjust the object transparency, and can toggle between overlaying the full rendering and overlaying just the edges (shown here).

## D  Controls

*3DB* takes an object-centric perspective, where an object of interest is spawned on a desired background. The scene mainly consists of the object and a camera. The controls in our pipeline affect this interplay between the scene components through various combinations of properties, which subsequently creates a wide variety of rendered images. The controls are implemented using the Blender Python API 'bpy' that exposes an easy to use framework for controlling Blender. 'bpy' primarily exposes a scene context variable, which contains references to the properties of the components such as objects and the camera; thus allowing for easy modification.

*3DB* comes with several predefined controls that are ready to use (see https://3db.github.io/3db/). Nevertheless, users are able (and encouraged) to implement custom controls for their use-cases.

# E   Additional Experiments Details

We refer the reader to our package `https://github.com/3db/3db` for all source code, 3D models, HDRIs, and config files used in the experiments of this paper.

For all experiments we used the pre-trained ImageNet ResNet-18 included in `torchvision`. In this section we will describe, for each experiment the specific 3D-models and environments used by *3DB* to generate the results.
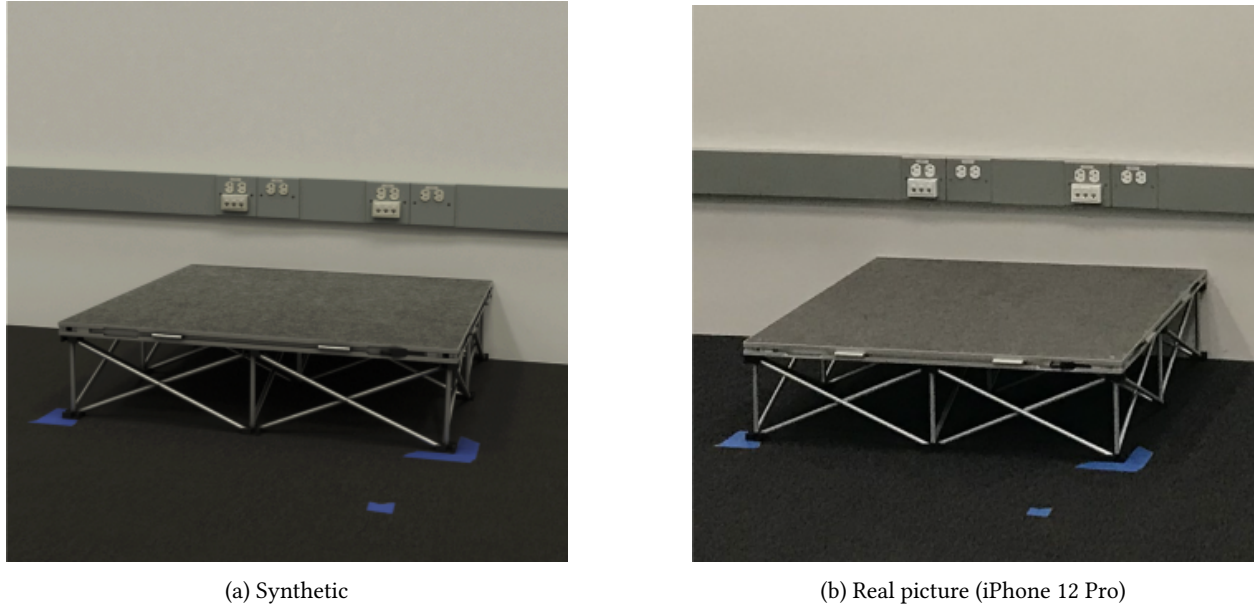


(a) Synthetic                                                                 (b) Real picture (iPhone 12 Pro)

Figure 20: Studio used for the real-world experiments (Section 4).

## E.1   Sensitivity to image backgrounds (Section 3.1)

### E.1.1   Analysing a subset of backgrdounds

**Models:**   We collected 19 3D-models in total. On top of the models shown on figure 23, we used models for: (1) an orange, (2) two different toy power drills, (3) a baseball ball, (4) a tennis ball, (5) a golf ball, (6) a running shoe, (7) a sandal and (8) a toy gun. Some of these models are from YCB [CWS+15] and the rest are purchased from `amazon.com` and then put through a 3D scanner to get corresponding meshes.

**Environments:**   We sourced 20 2k HDRI from the website `https://hdrihaven.com`. In particular we used: `abandoned_workshop`, `adams_place_bridge`, `altanka`, `aristea_wreck`, `bush_restaurant`, `cabin`, `derelict_overpass`, `dusseldorf_bridge`, `factory_yard`, `gray_pier`, `greenwich_park_03`, `kiara_7_late-afternoon`, `kloppenheim_06`, `rathaus`, `roofless_ruins`, `secluded_beach`, `small_hangar_02`, `stadium_01`, `studio_small_02`, `studio_small_04`.

### E.1.2   Analyzing all backgrounds with the "coffee mug" model.

**Models:**   We used a single model: the coffee mug, in order to keep computational resources under control.

**Environments:**   We used 408 HDRIs from `https://hdrihaven.com/` with a 2K resolution.

## E.2   Texture-shape bias (section 3.2)

**Textures:**   To replace the original materials, we collected 7 textures on the internet and we modified them to make them seamlessly tilable. These textures are shown on Figure 23.

**Models:**  We used all models that are shown on Figure 23.

**Environments:**  We used the virtual studio environment (Figure 20).

## E.3  Orientation and scale dependence (Section 3.3)

We use the same models and environments that are used in Appendix E.1.1.

## E.4  3D models Heatmaps (Figure 12)

**Models:**  For this experiment we used the set of models shown on Figure 23.

**Environments:**  We used the virtual studio environment (see Figure 20).

## E.5  Case study: using 3DB to dive deeper (Section 3.4)

**Models:**  We only used the mug since this experiment is mug specific.

**Environments:**  We used the sudio set shown on Figure 20.

## E.6  Physical realism (Section 4)

**Real-world pictures:**  All images were taken with an handheld Apple iPhone 12 Pro. To help us align the shots we used the application described in appendix C.

**Models:**  We used the models shown in Figure 15.

**Environments:**  The environment shown on Figure 20 was especially designed for this experiment. The goal was to have an environment that matches our studio as closely as possible. The geometry and materials were carefully reproduced using reference pictures. The lighting was reproduce through a high resolution HDRI map.

## E.7  Performance scaling (Appendix A)

The only relevant details for this experiment are the fact that we ran 10 policies (at most 5 concurrently). Each policy consisted of 1000 renders using a 2k HDRI as environment.
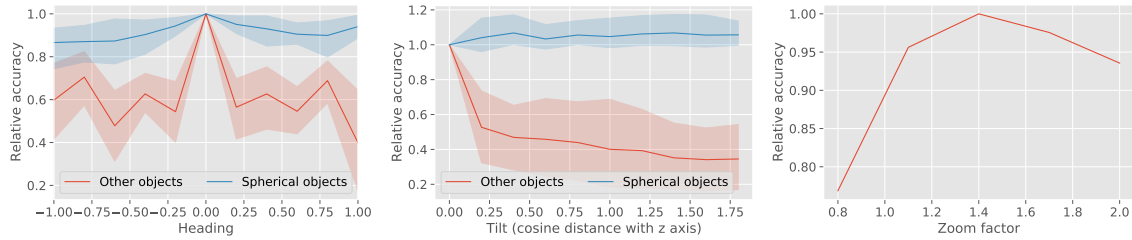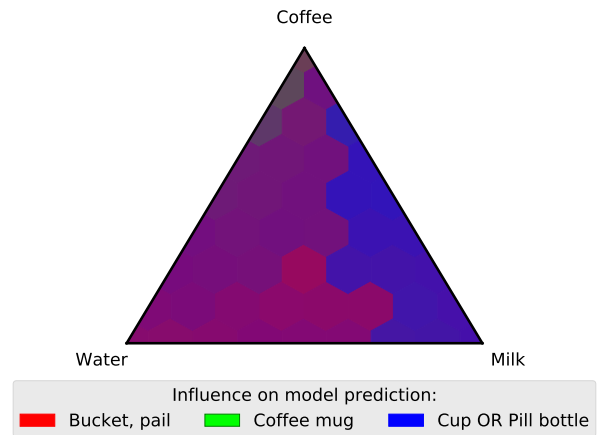
# F  Omitted Figures



Figure 21: Additonal plots to Figure 13. We plot the distribution of model accuracy as a function of object heading (*top*), tilt (*middle*) and zoom (*bottom*), aggregated over variations in controls. For heading and tilt, we separately evaluate accuracy for (non-)spherical objects. Notice how the performance of the model degrades for non-spherical objects as the heading/tilt changes, but not for spherical objects. Also notice how the performance depends on the zoom level of the camera (how large the object is in the frame).



(a) Sample of the images rendered for the experiment presented in section 3.4.

(b) Un-normalized version of Figure 14-(b).

Figure 22: Additional illustration for the mug liquid experiment of Figure 14. This figure shows the correlation of the liquid mixture in the mug on the prediction of the model, averaged over random viewpoints
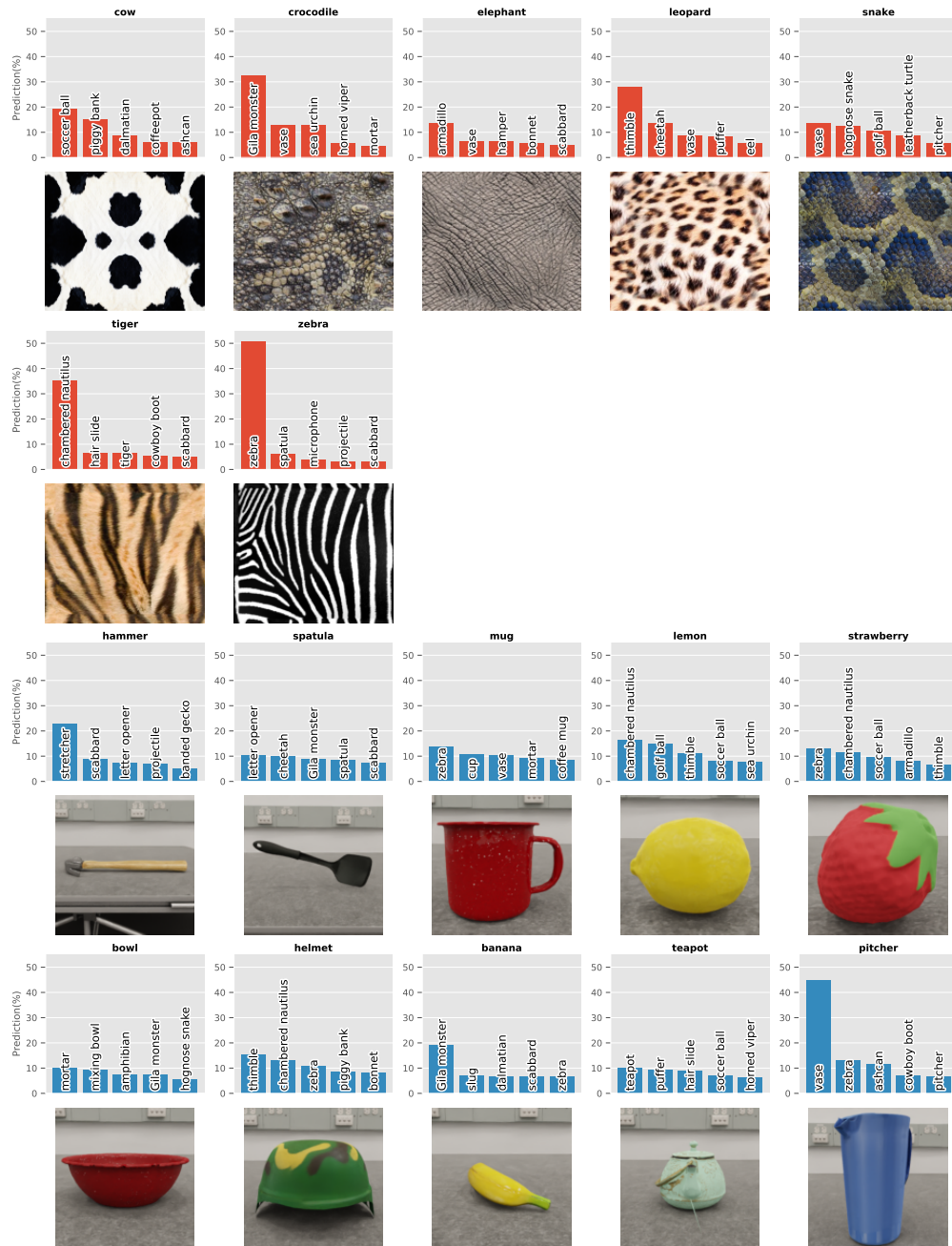
23

Figure 23: Additional examples of the experiment in Figure 11. Distribution of classifier predictions after the texture of the 3D object model is altered. In the top rows, we visualize the most frequently predicted classes for each texture (averaged over all objects). In the bottom rows, we visualize the most frequently predicted classes for each object (averaged over all textures). We find that the model tends to predict based on the texture more often than based on the object.