# Small in-distribution changes in 3D perspective and lighting fool both CNNs and Transformers

**Spandan Madan**
School of Engineering and Applied Sciences
Harvard University
Cambridge, MA 02138, USA
spandan_madan@seas.harvard.edu

**Tomotake Sasaki**
Artificial Intelligence Laboratory
Fujitsu Limited
Kawasaki, Kanagawa 211-8588, Japan
tomotake.sasaki@fujitsu.com

**Tzu-Mao Li** *
Computer Science and Engineering
University of California, San Diego
San Diego, CA 92093, USA
tzli@ucsd.edu

**Xavier Boix** *
Department of Brain and Cognitive Sciences
Massachusetts Institute of Technology
Cambridge, MA 02139, USA
xboix@mit.edu

**Hanspeter Pfister**
School of Engineering and Applied Sciences
Harvard University
Cambridge, MA 02138, USA
pfister@seas.harvard.edu

## Abstract

Neural networks are susceptible to small transformations including 2D rotations and shifts, image crops, and even changes in object colors. This is often attributed to biases in the training dataset, and the lack of 2D shift-invariance due to not respecting the sampling theorem. In this paper, we challenge this hypothesis by training and testing on unbiased datasets, and showing that networks are brittle to both small 3D perspective changes and lighting variations which cannot be explained by dataset bias or lack of shift-invariance. To find these in-distribution errors, we introduce an evolution strategies (ES) based approach, which we call *CMA-Search*. Despite training with a large-scale ($\sim 0.5$ million images), unbiased dataset of camera and light variations, in over 71% cases *CMA-Search* can find camera parameters in the vicinity of a correctly classified image which lead to in-distribution misclassifications with $< 3.6\%$ change in parameters. With lighting changes, CMA-Search finds misclassifications in 33% cases with $< 11.6\%$ change in parameters. Finally, we extend this method to find misclassifications in the vicinity of ImageNet images for both ResNet and OpenAI's CLIP model.

## 1   Introduction

Neural networks models are highly susceptible to seemingly benign changes—two dimensional rotations and translations [1], image crops [2, 3], and even changes in the color space [4, 5, 6].
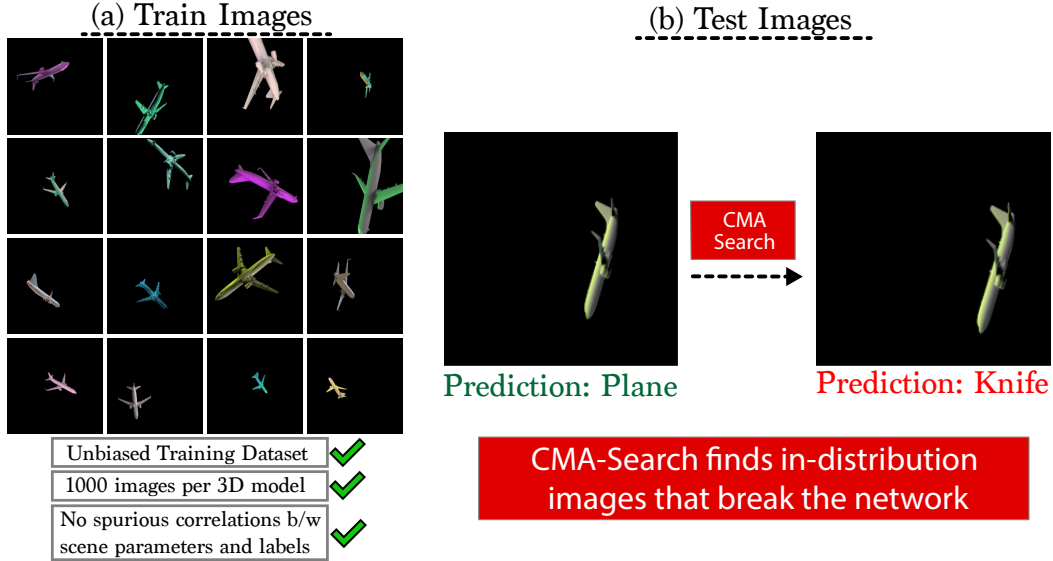
---

*equal advising contribution

Figure 1: *In-distribution failures discovered by CMA-Search despite unbiased training data.* (a) Sample images from our large scale and unbiased dataset of 3D objects seen under camera and lighting variations. There are 11 categories with 40 3D models each and 1000 images per 3D model resulting in $\sim 0.5$ million images. (b) Our gradient-free, evolution strategies based approach finds failures by searching the vicinity of camera parameters. Unlike typical adversarial attacks, our method does not add noise, and our constraints ensure that identified errors are in-distribution.

This is concerning as these transformations occur naturally in the real world, unlike the synthetic perturbations typically added in adversarial attacks. These susceptibilities are often attributed to the dataset bias (systematic differences between the training and testing distributions) in ImageNet and CIFAR [7, 8], and to the lack of shift invariance due to not respecting the sampling theorem [9, 10, 3]. Here, we question if this hypothesis is indeed the complete picture, or if there is more underlying complexity which remains unexplored. Concretely, we train and test visual recognition models across variations in camera and lighting, while ensuring large-scale, unbiased training data with no spurious correlations between the camera, lighting and the image labels. To investigate the role of shift-invariance, multiple shift-invariant architectures are also trained.

These experiments are enabled by our computer graphics pipeline for generating and modifying images by directly accessing the scene's camera and lighting parameters. We use our pipeline to create a large-scale ($\sim 0.5$ million images) and unbiased dataset of rendered ShapeNet [11] objects with camera and lighting variations—Fig. 1(a) shows examples of the same 3D model under different conditions, while Fig. 2(b) shows 3D models for 4 categories. This dataset contains objects seen from multiple viewpoints, shifted across the frame, and we use physically based rendering [12, 13] to accurately simulate complex lighting artifacts including multiple colors and self-shadows which makes the dataset challenging for neural networks as corroborated in Table. 1.

Investigating the brittleness of networks requires densely searching the space of camera and lighting parameters for failure cases. We apply an evolution-strategies based approach to propose a new search for in-distribution, misclassified images which we call *CMA-Search*. Starting with a correctly classified image, our method looks in the vicinity of the camera and lighting parameters to find an in-distribution image which is incorrectly classified. Note that unlike adversarial attacks, our method does not add noise and our constraints ensure that identified errors are in-distribution.

Analyzing the brittleness using *CMA-Search*, we find that visual recognition models are extremely brittle to both changes in 3D perspective due to camera movement and lighting. Our method can find misclassified images in the vicinity of correctly classified images in over 71% cases, with less than 3.6% change in the camera position. With lighting changes, *CMA-Search* can find a

misclassification in 33% cases with $< 5\%$ change in light position. Furthermore, by combining our search method with a recent single-view view synthesis model [14], we show similar results of brittleness with ImageNet images for both a pretrained ResNet [15], and the recent start-of-the-art transformer based OpenAI CLIP architecture [16]. By comparing against 2D shifts, we show that while shift-invariant architectures are adept at handling 2D shifts, they are still very susceptible to small 3D perspective changes. These results are consistent across multiple state-of-the-art CNNs and transformer architectures. In fact, we report that with a dataset of $\sim 0.5$ million images, transformers struggle much more than CNNs at performing well across variations in camera and lighting. The code to reproduce these results, the data and the computer graphics pipeline will be made available upon publication.

## 2 Related Work

### 2.1 Susceptibility of neural networks to small 2D transformations

Susceptibility to small transformations [3], crops [2], and 2D rotations and translations [1] are often attributed to lack of shift invariance in modern CNNs [7, 17, 18, 19]. Thus, proposed alternative architectures have focused on being shift invariant—anti-aliasing networks use the seminal signal processing trick of anti-aliasing [9], while recently proposed truly shift invariant networks propose a new sampling methodology to guarantee a 100% consistency in classification under 2D shifts [10]. Unlike our work, these works have focused on only 2D transformations.

### 2.2 Semantic adversarial attacks and in-distribution brittleness

To understand the failure modes of neural networks described in Sec. 2.1, recent work has sought to generate adversarial perturbations which are human interpretable i.e. semantic adversarial examples. These works often rely on synthetic data, using differentiable rendering or other optimization methods to find adversarial images by modifying scene parameters [5, 20, 21, 22, 23, 24, 25, 26]. These include a custom differentiable renderer to perturb the camera, lighting, or object mesh vertices [20], and using a neural renderer where light is represented by network activations [21].

While we rely on computer graphics as well, there is a key differences between these works and our paper—we ensure that our train and test distributions match, so all our identified errors are in-distribution. Above works do not attempt to search for in-distribution errors and add noise which is not constrained to be in-distribution. In contrast, our approach (*CMA-Search*) searches only within the space of camera and lighting parameters to find in-distribution errors without adding any noise.

### 2.3 Applications of computer graphics beyond adversarial attacks

Several recent works have turned to computer graphics (CG) as a means to generate synthetic datasets which may be hard to create in the real world. Originally, synthetic data was primarily used to increase the size of training data by data augmentation [27, 28, 29, 30, 31]. However, recent years have seen a shift in this trend, with several works using computer graphics to create controlled environments for training, testing and studying the generalization capabilities of neural networks [27, 32, 33, 34, 35, 36]. Recent work has also shown that behaviour identified on controlled, synthetic datasets extends well to natural data [33, 28, 29]. Inspired by this, we use a computer graphics pipeline to generate controlled, unbiased data with camera and light variations.

## 3 Generating an unbiased training dataset of camera and light variations

To understand the role of dataset bias on the brittleness of CNNs to viewpoints, rotations and translations, or color changes, many works have suggested the use of large-scale, unbiased datasets [7, 3]. However, collecting such a dataset presents many difficulties. Firstly, as the same object must be seen under multiple camera and lighting variations, such a dataset cannot be created by scraping the internet—these objects must be photographed. However, photographing enough objects under varied viewpoint and lighting conditions to train data-hungry models like vision transformers which require millions of images is challenging. Secondly, to ensure there are no spurious correlations between object category and nuisance factors like textures, viewpoints, and lighting in the real
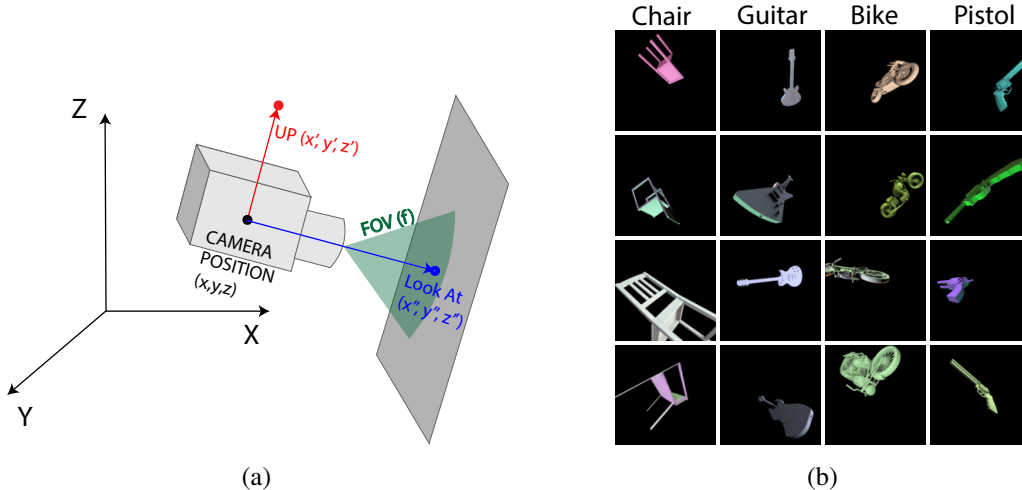
|   | Chair | Guitar | Bike | Pistol |
|---|---|---|---|---|

(a)                                             (b)

Figure 2: *3D scene setup and resulting images*. (a) Images in our dataset are completely parametrized by the camera and light. Physical interpretation of the camera parameters is illustrated here. Analogously, light is parametrized by the position, look at, 2D size and the RGB intensities. (b) Sample images for 4 object categories generated using our 3D scene setup. As can be seen, images contain complex viewpoints and locations, multiple colors per object and complex artifacts like self-shadows.

world presents another serious challenge. So, inspired by the growing trend to use synthetic data for investigating generalization capabilities of neural networks, we prevent these issues by turning to computer graphics [27, 32, 33, 34, 35, 36]. Our graphics pipeline (explained below) easily allows us to generate a large-scale, unbiased dataset of objects seen under varied camera and lighting conditions.

**3D Scene Setup**: Each scene contains one camera, one 3D model and 1-4 lights. To ensure no spurious correlations with object texture [17], texture for all ShapeNet objects was replaced with a simple diffuse material and the background was kept constant to ensure no spurious correlations between foreground and background. Thus, every scene is completely parametrized by the camera and the light parameters. As shown in Fig. 2(a), camera parameters are 10 Dimensional: one dimension for the **FOV** (field of view of camera lens), and three dimensions each for the **Camera Position** (coordinates of camera center), **Look At** (point on the canvas where the camera looks), and the **UP vector** (rotation of camera). Analogously, lights are represented by 11 dimensions - two dimensions for the Light Size, and three each for Light Position, light Look At and RGB color intensity. Multiple lights ensure that scenes contain complex mixed lighting, including self-shadows. Thus, our scenes are $(11n + 10)$ dimensional, where $n$ is the number of lights. There is a one-to-one mapping between the pixel space (rendered images) and this low dimensional scene representation.

**Unbiased, uniformly distributed training dataset:** Our dataset contains 11 categories, with 40 3D models for every category chosen from ShapeNet [11]. Sample images from four categories are shown in Fig. 2(b). Each 3D model was rendered under 1000 different camera and lighting conditions following the scene setup described above. Fig. 1(a) shows one 3D model rendered under different conditions. To ensure a good distribution over viewpoints, locations, perspective projections and colors, it was ensured that scene parameters follow a uniform distribution. Concretely, camera and light positions were sampled from a uniform distribution on a spherical shell with a fixed minimum and maximum radius. The Up Vector was uniformly distributed across range of all possible camera rotations, and RGB light intensities were uniformly distributed across all possible colors. Camera and light Look At positions were uniformly distributed while ensuring the object stays in frame and is well-lit (frame size depends on Camera Position and FOV). Finally, Light Size and camera FOV were uniformly sampled 2D and 1D vectors with pre-defined upper and lower bounds. More details on implementing uniform sampling in these domains are provided in the supplement. All networks were trained with a learning rate of 0.0003, with an Adam optimizer using PyTorch and NVIDIA TeslaK80 GPUs on a compute cluster. More details are provided in the supplement.

4

**Challenging test data to evaluate in-distribution performance:** Networks trained on the unbiased dataset described above are evaluated on two test sets - one with the 3D models seen during training, and the second with new, unseen 3D models. The first test set was generated by simply repeating the same procedure as described in Sec. 3. Thus, the *(Geometry × Camera × Lighting)* joint distribution matches exactly for the train set and this test set. The second test set was created by the exact same generation procedure, but with 10 new 3D models for every category chosen from ShapeNet [11]. The motivation for this second test set was two-fold: Firstly, evaluating generalization performance while ensuring that all nuisance parameters (viewpoint, lighting) are exactly matched in training and testing. Secondly, to ensure our models are not over-fitting to the 3D models used for training.

# 4   CMA-Search: Finding in-distribution failures by searching the vicinity

To investigate the brittleness of neural networks with respect to changes in camera and lighting, we propose a new, gradient-free search method to find incorrectly classified images. Starting with a correctly classified image, our method searches the vicinity by slightly modifying camera or light parameters to find an in-distribution error. While adversarial viewpoints and lighting have been reported before in the literature [20, 21, 23], there are two major differences in our approach. First, these methods search for an adversarial image by adding noise to the image without constraining the resulting image to be within the training distribution. In comparison, our approach does not add noise, but instead searches within the distribution to find in-distribution errors. Secondly, unlike our gradient-free search method, these methods often rely on gradient descent and thus require high dimensional representations of the scene to work well. For instance, these works often use neural rendering where network activations act as a high dimensional representation of the scene [21, 25], or use up-sampling of meshes to increase dimensionality [20].

To extend these approaches to work well with our low-dimensional scene representation, we apply a gradient-free optimization method to search the space—Covariance Matrix Adaptation-Evolution Strategy (CMA-ES) [37, 38]. We found that gradient descent with differentiable rendering struggled to find in-distribution errors in our scenes due to the low dimensionality of the optimization problem. CMA-ES has been found to work reliably well with non-smooth optimization problems and especially with local optimization [39], which made it a perfect fit for our search strategy. In Fig. 3 we show examples of in-distribution failures found by our *CMA-Search* method. Starting with the correctly classified image (left), our method finds an image in the vicinity by slightly modifying camera parameters of the scene. As can be seen, subtle changes in 3D can lead to drastic errors in classification. We also highlight the subtle changes in camera position (in black) and camera Look At (in blue) in the figure. To the best of our knowledge, this is the first evolutionary strategies based search method for finding in-distribution failures, and it also works well in low dimensions.

Starting from the initial scene parameters, CMA-ES generates offspring by sampling from a multivariate normal (MVN) distribution i.e. mutating the original parameters. These offspring are then sorted based on the fitness function (classification probability), and the best ones are used to modify the mean and covariance matrix of the MVN for the next generation. The mean represents the current best estimate of the solution i.e. the maximum likelihood solution, while the covariance matrix dictates the direction in which the population should be directed in the next generation. The search is stopped either when a misclassification occurs, or after 15 iterations. Algorithm 1 provides an outline for the method which was implemented using pycma [40]. The algorithm for searching light parameters which lead to misclassification is analogous. More details on the parameter update subroutines are provided in the supplement.

# 5   Results

We present results on the distribution of classification errors made by visual recognition models as camera and lights are varied. To characterize this distribution in detail, we report both the accuracy on randomly sampled images from the test distribution, and the brittleness of the networks using *CMA-Search* which searches the vicinity of correctly classified images to find failures. For the former, we evaluate these networks on an ImageNet sized test sample (40,000 images) covering the whole space of camera and light variations.
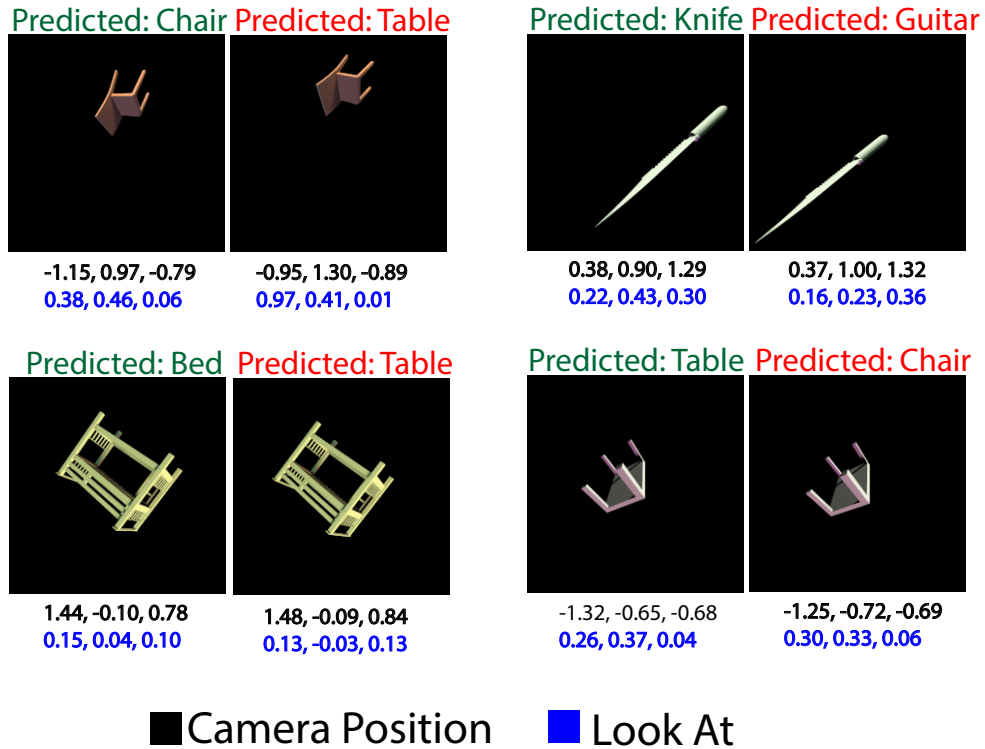
| Predicted: Chair | Predicted: Table | | Predicted: Knife | Predicted: Guitar |

-1.15, 0.97, -0.79     -0.95, 1.30, -0.89          0.38, 0.90, 1.29      0.37, 1.00, 1.32
0.38, 0.46, 0.06       0.97, 0.41, 0.01            0.22, 0.43, 0.30      0.16, 0.23, 0.36

| Predicted: Bed | Predicted: Table | | Predicted: Table | Predicted: Chair |

1.44, -0.10, 0.78      1.48, -0.09, 0.84          -1.32, -0.65, -0.68   -1.25, -0.72, -0.69
0.15, 0.04, 0.10       0.13, -0.03, 0.13           0.26, 0.37, 0.04      0.30, 0.33, 0.06

■ Camera Position     ■ Look At

Figure 3: *CMA-Search across camera parameters.* Starting with the correctly predicted images, our evolutionary strategies based method searches the vicinity of camera parameters for subtle 3D perspective changes that lead to misclassification. Note that these are in-distribution failures, and not adversarial examples which are often generated by adding noise without in-distribution constraints. We find that networks are most sensitive to changes in camera position and look at parameters, and the subtle changes in these parameters are reported here alongside misclassified images.

## 5.1 Networks struggle to generalize across camera and light variations

Table 1 reports accuracy for several state-of-the-art CNNs [15, 9, 10] and transformer architectures including the vision transformer (ViT) [41], and the data efficient transformer (DeIT) and its distilled version (DeIT Distilled) [42]. As can be seen, there is still much room for improvement. That is, neural networks do not perform well across camera and lighting variations despite ensuring: (1) uniformly distributed and unbiased training data, (2) 1000 images per 3D object (total 0.5 million images), and (3) no spurious correlations between the scene parameters and the image labels. In fact, the problem is even more pronounced with transformer models. To test if this problem can be mitigated with shift-invariant architectures, we also report results on two specialized shift-invariant architectures - Anti-Aliased Networks [9], and the recent Truly Shift Invariant Network [10]. While these networks do provide a boost in performance, they too are susceptible to camera and lighting variations. Furthermore, we find that our neural networks have not overfit and generalize well to new 3D models. However, the performance on these new 3D models also mirrors the same trend.

These results naturally raise the question—What images are these networks failing on? Are there certain lighting and camera conditions that the networks fail on? The one-to-one mapping between the pixel space (images) and our low-dimensional scene representation allows us to answer these questions by visualizing and comparing correctly and incorrectly classified images in this low dimensional space. In Fig. 4 we show the distribution of camera and lighting parameters for images which were classified incorrectly. As can be seen, the errors seem well distributed across space—we found no clear, strong patterns which characterize the camera and light conditions of misclassified images.

6

**Algorithm 1** CMA-ES based camera parameter search for in-distribution failures.

1: Let $x \in \mathbb{R}^{10}$ denote camera parameters.
2: **function** FITNESS($x$)
3:     image = *Render*($x$)
4:     predicted_category, probability = *Network*(image)
5:     **return** predicted_category, probability

6:
7: Let $x_{init}$ denote initial camera parameters, $\lambda$ be number offspring per generation, and $y$ be the image category.
8:
9: **procedure** CMA-SEARCH($x_{init}, \lambda, y$)
10:     **initialize** $\mu = x_{init}, C = I$             $\triangleright$ $I$ denotes identity matrix.
11:     **while** True **do**
12:         **for** j in $\lambda$ **do**
13:             $x_j$ = sample_multivariate_normal($\mu, C$)     $\triangleright$ Generate mutated offspring
14:             $y_j, p_j$ =FITNESS($x_j, R, N$)     $\triangleright$ Calculate fitness of offspring
15:             **if** $y_j \neq y$ **then**
16:                 **return** $x'$     $\triangleright$ Classification fails for image with camera parameters $x'$
17:         $x_{1...\lambda} \leftarrow x_{s(1)...s(\lambda)}$, with $s(j)$ = argsort($p_j$)     $\triangleright$ Pick best offspring
18:         $\mu, C \leftarrow$ update_parameters($x_{1...\lambda}, \mu, C$)

Table 1: Performance of visual recognition models on seen and new 3D models.

| Accuracy | ResNet | ResNet (pretrained) | Anti-Aliased Networks | Truly Shift Invariant | ViT | DeIT | DeIT Distilled |
|---|---|---|---|---|---|---|---|
| Seen models | 0.75 | 0.76 | 0.82 | 0.80 | 0.58 | 0.63 | 0.64 |
| New models | 0.70 | 0.70 | 0.74 | 0.72 | 0.59 | 0.64 | 0.65 |



Camera Positions (Incorrect Predictions)    Look At (Incorrect Predictions)    Up Vector (Incorrect Predictions)    Field of View (Incorrect Predictions)
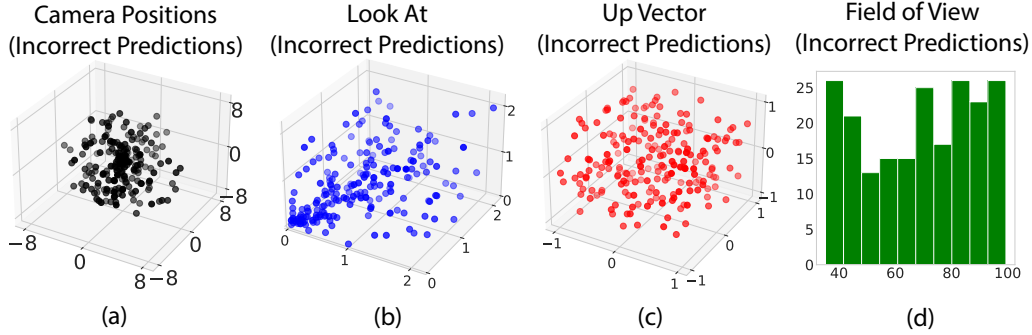
(a)      (b)      (c)      (d)

Figure 4: *Errors are well distributed across the scene parameter space* (a) Coordinates of camera positions, (b) Coordinates of Look At, (c) Up Vector and (d) Histogram of errors across lens field of view. We found no clear, strong patterns which characterize the camera and light conditions of misclassified images. This is in contrast to human vision which is impacted by regions of camera positions (non-canonical viewpoints), and up vector (upside-down orientations) among others.

Note that regions in each of these parametric spaces represent human interpretable scenarios which have been known to impact human vision significantly. For instance, changes in camera position represent canonical vs non-canonical poses which significantly impact human vision [43, 44, 45]. Similarly, changes in the up vector can represent upside-down objects which too impact human vision [46, 47, 48]. In contrast, Fig. 4 shows that networks do not suffer in specific regions of the space. These observations corroborate our finding from *CMA-Search*: correctly classified and misclassified images are in the vicinity of each other in the space of camera and light parameters. This result is consistent across multiple architectures (see supplement).

Table 2: Accuracy under 2D shifts and with *CMA-Search.*

| Accuracy | ResNet | ResNet (pretrained) | Anti-Aliased Networks | Truly Shift Invariant | ViT | DeIT | DeIT Distilled |
|----------|--------|---------------------|-----------------------|-----------------------|-----|------|----------------|
| 2D shifts | 0.82 | 0.80 | 0.90 | 0.96 | 0.47 | 0.47 | 0.59 |
| **CMA Cam** | **0.29** | **0.42** | **0.55** | **0.47** | **0.15** | **0.15** | **0.14** |
| **CMA Light** | **0.67** | **0.77** | **0.80** | **0.76** | **0.46** | **0.59** | **0.48** |

## 5.2 Brittleness: networks are susceptible to small changes in 3D perspective and lighting

As shown in Table 2, *CMA-Search* finds small changes in 3D perspective and lighting which have a drastic impact on network performance. Starting with an image correctly classified by a ResNet18 model, our method can find an error in its vicinity for 71% cases with $< 3.6\%$ change in the camera position, and $< 11.6\%$ change in the camera Look At. For transformers, the impact is far worse. There were only about 15% images for which *CMA-Search* could not find an error in the vicinity. Similarly, with lighting changes *CMA-Search* can find a misclassification in 33% cases with $< 5\%$ change in light position as shown in Table 2. In the supplement we provide additional figures reporting the percent of change required in camera and light parameters to break these architectures.

Investigating these errors, we find that networks are most sensitive to changes in the Camera Position and the camera Look At—small, in-distribution 3D perspective changes. Comparing these transformations against 2D shifts, we find that the impact of 2D shifts on a ResNet model is not as drastic. In fact, shift-invariant architectures are largely unaffected by 2D shifts, with no errors found in the vicinity for almost 95% cases. However, these architectures too are highly susceptible to both 3D perspective changes, and subtle lighting changes.

These results reveal that the space of camera and lighting variations is filled with in-distribution failures. Recall that previous works have suggested that the brittleness of neural networks stems from biased data or the lack of shift invariance. However, our results here show that this in-distribution brittleness occurs despite training shift-invariant architectures with an unbiased dataset. Thus, an unbiased dataset and shift-invariant architectures are necessary, but not sufficient conditions to build systems which can gracefully handle camera and lighting variations.

## 6   Results on ImageNet with ResNet and OpenAI CLIP networks

So far we have focused on images generated by our graphics pipeline as extending these results to natural image datasets presents a challenge—generating images in the vicinity of a correctly classified image by slightly modifying the camera parameters. To do so for ImageNet is equivalent to novel view synthesis (NVS) from single images. This is a highly challenging task but recent advances in NVS enable us to extend our method to natural image datasets like ImageNet [49, 50, 51, 14].

To generate new views in the vicinity of ImageNet images, we rely on a single-view synthesis model based on multi-plane images (MPI) [14]. The MPI model takes as input an image and the $(x, y, z)$ offsets which describe camera movement along the X, Y and Z axes. Note that unlike our renderer, it cannot introduce changes to the camera Look At, Up Vector, Field of View or lighting changes. As before, we used *CMA-Search* to optimize the camera parameters, but instead of our renderer, we now use the view synthesis model for generating novel views of ImageNet images. The MPI model was not trained on ImageNet, and can at times fail to generate novel views, resulting in blurry images instead. We omit these images to only present results on failures due to small, 3D perspective changes. Starting with a correctly predicted ImageNet image, we use *CMA-Search* in conjunction with the MPI model to find images in the vicinity with small, 3D perspective changes which can break ImageNet trained classification networks including ResNet18, and OpenAI's transformer based CLIP model [16]. Results for these experiments are reported in Fig. 5 with more examples provided in the supplement. While these results present an interesting application on natural images, note that we cannot be entirely sure that the images found by *CMA-Search* on imagenet are indeed in-distribution, further proving the utility of our computer graphics based approach.
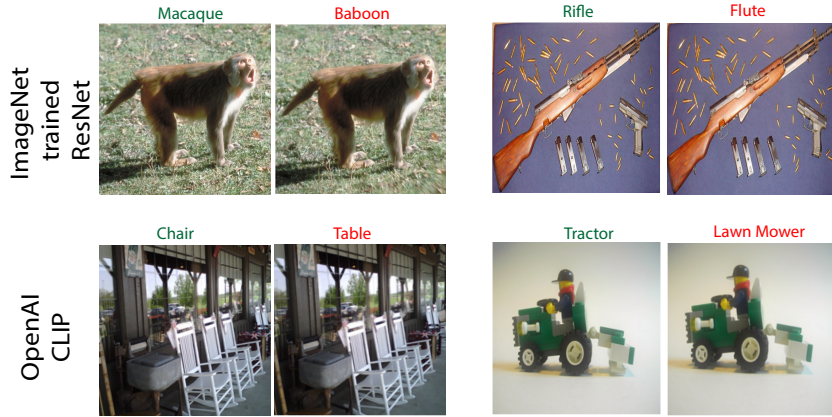
Figure 5: *CMA-Search on ImageNet images.* To replicate results on ImageNet, we replace our rendering pipeline with the single view MPI [14] model to generate novel views of ImageNet images. Here we show results using *CMA-Search* with the MPI model to find subtle 3D perspective changes which lead to misclassification with ResNet18 and OpenAI's CLIP model.

## 7 Conclusions

Susceptibilities of recognition models have often been attributed to biased training data and the lack of shift invariance in modern CNNs. We put this hypothesis to test by training and testing with a large-scale, unbiased dataset and propose a new search method for investigating the brittleness of neural networks. Our findings show that while data augmentation, unbiased datasets, and specialized shift-invariant architectures would certainly be helpful, the real problem runs far deeper—networks are susceptible to small changes in 3D perspective and lighting. The real world is rife with these transformations. Thus, building on 2D shift-invariance literature, our work suggests moving towards a new class of architectures which are also invariant to a super-set of 2D shifts—3D perspective changes and lighting. A promising thread of works has focused on the closely-related problem of 3D viewpoint invariance [52, 53, 54, 55]. We believe that incorporating these ideas into modern architectures which scale well to large datasets like ImageNet [56, 57] and Google's JFT-300M [58] opens a promising new direction which can help make networks more robust to these transformations.

## References

[1] Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. A rotation and a translation suffice: Fooling CNNs with simple transformations. `https://openreview.net/forum?id=BJfvknCqFQ`, 2018.

[2] Sanjana Srivastava, Guy Ben-Yosef, and Xavier Boix. Minimal images in deep neural networks: Fragile object recognition in natural images. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.

[3] Aharon Azulay and Yair Weiss. Why do deep convolutional networks generalize so poorly to small image transformations? *Journal of Machine Learning Research*, 20(184):1–25, 2019.

[4] Jeet Mohapatra, Tsui-Wei Weng, Pin-Yu Chen, Sijia Liu, and Luca Daniel. Towards verifying robustness of neural networks against a family of semantic perturbations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 244–252, 2020.

[5] Hossein Hosseini and Radha Poovendran. Semantic adversarial examples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1614–1619, 2018.

[6] Ali Shahin Shamsabadi, Ricardo Sanchez-Matilla, and Andrea Cavallaro. Colorfool: Semantic adversarial colorization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1151–1160, 2020.

[7] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. In *Advances in Neural Information Processing Systems*, 2019.

[8] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1521–1528. IEEE, 2011.

[9] Richard Zhang. Making convolutional networks shift-invariant again. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 7324–7334, 2019.

[10] Anadi Chaman and Ivan Dokmanić. Truly shift-invariant convolutional neural networks. *arXiv preprint arXiv:2011.14214*, 2020.

[11] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. ShapeNet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015.

[12] Tzu-Mao Li, Miika Aittala, Frédo Durand, and Jaakko Lehtinen. Differentiable Monte Carlo ray tracing through edge sampling. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018.

[13] Matt Pharr, Wenzel Jakob, and Greg Humphreys. *Physically based rendering: From theory to implementation*. Morgan Kaufmann, 2016.

[14] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 551–560, 2020.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.

[16] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.

[17] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.

[18] Jason Jo and Yoshua Bengio. Measuring the tendency of CNNs to learn surface statistical regularities. *arXiv preprint arXiv:1711.11561*, 2017.

[19] Taro Makino, Stanislaw Jastrzebski, Witold Oleszkiewicz, Celin Chacko, Robin Ehrenpreis, Naziya Samreen, Chloe Chhor, Eric Kim, Jiyon Lee, Kristine Pysarenko, et al. Differences between human and machine perception in medical diagnosis. *arXiv preprint arXiv:2011.14036*, 2020.

[20] Hsueh-Ti Derek Liu, Michael Tao, Chun-Liang Li, Derek Nowrouzezahrai, and Alec Jacobson. Beyond pixel norm-balls: Parametric adversaries using an analytically differentiable renderer. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.

[21] Xiaohui Zeng, Chenxi Liu, Yu-Siang Wang, Weichao Qiu, Lingxi Xie, Yu-Wing Tai, Chi-Keung Tang, and Alan L Yuille. Adversarial attacks beyond the image space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4302–4311, 2019.

[22] Rakshith Shetty, Mario Fritz, and Bernt Schiele. Towards automated testing and robustification by semantic adversarial data generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 489–506, 2020.

[23] Lakshya Jain, Steven Chen, Wilson Wu, Uyeong Jang, Varun Chandrasekaran, Sanjit Seshia, and Somesh Jha. Generating semantic adversarial examples with differentiable rendering. `https://openreview.net/forum?id=SJlRF04YwB`, 2019.

[24] Chaowei Xiao, Dawei Yang, Bo Li, Jia Deng, and Mingyan Liu. Meshadv: Adversarial meshes for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6898–6907, 2019.

[25] Ameya Joshi, Amitangshu Mukherjee, Soumik Sarkar, and Chinmay Hegde. Semantic adversarial attacks: Parametric transformations that fool deep classifiers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4773–4783, 2019.

[26] Philip Yao, Andrew So, Tingting Chen, and Hao Ji. On multiview robustness of 3D adversarial attacks. In *Practice and Experience in Advanced Research Computing*, pages 372–378. 2020.

[27] Weichao Qiu and Alan Yuille. UnrealCV: Connecting computer vision to Unreal Engine. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 909–916, 2016.

[28] Gül Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning from synthetic humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[29] Jiteng Mu, Weichao Qiu, Gregory D Hager, and Alan L Yuille. Learning from synthetic animals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12386–12395, 2020.

[30] Michael A Alcorn, Qi Li, Zhitao Gong, Chengfei Wang, Long Mai, Wei-Shinn Ku, and Anh Nguyen. Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4845–4854, 2019.

[31] Matthew Johnson-Roberson, Charles Barto, Rounak Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram Vasudevan. Driving in the Matrix: Can virtual worlds replace human-generated annotations for real world tasks? In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2017.

[32] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.

[33] Spandan Madan, Timothy Henry, Jamell Dozier, Helen Ho, Nishchal Bhandari, Tomotake Sasaki, Frédo Durand, Hanspeter Pfister, and Xavier Boix. On the capability of neural networks to generalize to unseen category-pose combinations. *arXiv preprint arXiv:2007.08032*, 2020.

[34] Philipp Bomatter, Mengmi Zhang, Dimitar Karev, Spandan Madan, Claire Tseng, and Gabriel Kreiman. When pigs fly: Contextual reasoning in synthetic and natural scenes. *arXiv preprint arXiv:2104.02215*, 2021.

[35] Guillaume Leclerc, Hadi Salman, Andrew Ilyas, Sai Vemprala, Ashish Kapoor, and Aleksander Madry. 3DB: A framework for analyzing computer vision models with simulated data. `https://github.com/3db/3db`, 2021.

[36] Jingkang Wang, Ava Pun, James Tu, Sivabalan Manivasagam, Abbas Sadat, Sergio Casas, Mengye Ren, and Raquel Urtasun. AdvSim: Generating safety-critical scenarios for self-driving vehicles. *arXiv preprint arXiv:2101.06549*, 2021.

[37] Nikolaus Hansen and Andreas Ostermeier. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation. In *Proceedings of IEEE International Conference on Evolutionary Computation*, pages 312–317, 1996.

[38] Nikolaus Hansen. The CMA evolution strategy: A tutorial. *arXiv preprint arXiv:1604.00772*, 2016.

[39] Nikolaus Hansen and Andreas Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9(2):159–195, 2001.

[40] Nikolaus Hansen, Youhei Akimoto, and Petr Baudis. CMA-ES/pycma on Github. Zenodo, DOI:10.5281/zenodo.2559634, February 2019.

[41] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.

[42] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020.

[43] Pablo Gomez, Jennifer Shutter, and Jeffrey N Rouder. Memory for objects in canonical and noncanonical viewpoints. *Psychonomic Bulletin & Review*, 15(5):940–944, 2008.

[44] Kyla P Terhune, Grant T Liu, Edward J Modestino, Atsushi Miki, Kevin N Sheth, Chia-Shang J Liu, Gabrielle R Bonhomme, and John C Haselgrove. Recognition of objects in non-canonical views: A functional MRI study. *Journal of Neuro-Ophthalmology*, 25(4):273–279, 2005.

[45] Volker Blanz, Michael J Tarr, and Heinrich H Bülthoff. What object attributes determine canonical views? *Perception*, 28(5):575–599, 1999.

[46] Wolfgang Köhler. *Dynamics in psychology*. WW Norton & Company, 1960.

[47] Michael B Lewis. The lady's not for turning: Rotation of the Thatcher illusion. *Perception*, 30(6):769–774, 2001.

[48] Peter Thompson. Margaret Thatcher: a new illusion. *Perception*, 1980.

[49] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5336–5345, 2020.

[50] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 286–301, 2016.

[51] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7467–7477, 2020.

[52] Robert Gens and Pedro M Domingos. Deep symmetry networks. *Advances in Neural Information Processing Systems*, pages 2537–2545, 2014.

[53] Fabio Anselmi, Lorenzo Rosasco, and Tomaso Poggio. On invariance and selectivity in representation learning. *Information and Inference: A Journal of the IMA*, 5(2):134–158, 2016.

[54] Qianli Liao, Joel Z Leibo, and Tomaso Poggio. Learning invariant representations and applications to face verification. In *Advances in Neural Information Processing Systems*, 2013.

[55] Tomaso Poggio and Fabio Anselmi. *Visual cortex and deep networks: learning invariant representations*. MIT Press, 2016.

[56] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.

[57] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.

[58] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 843–852, 2017.

# Supplementary Material

## A Uniform sampling of scene parameters for generating unbiased data

As explained in Sec.3 of the main text, to ensure that our dataset equally represents the whole range of object viewpoints, locations, colors, and other camera/light parameters, each camera and light parameter was sampled from a uniform distribution on a predefined, reasonable support. To support differentiable rendering and compare it with our CMA-ES based search method, we used the Redner to render our images [1]. Below we specify the hyper-parameters for rendering, along with how each of the parameters is sampled in our implementation.

**Camera Position:** For scene camera, first a random radius $r_c$ is sampled while ensuring $r_c \sim \mathrm{Unif}(0.5, 8)$. Then, the camera is placed on a random point denoted $(x_c, y_c, z_c)$ on the spherical shell of radius $r_c$. To generate a random point on the sphere while ensuring an equal probability of all points, we rely on the method which sums three randomly sampled normal distributions:

$$X, Y, Z \sim \mathcal{N}(0, 1),$$
$$v = (X, Y, Z),$$
$$(x_c, y_c, z_c) = r_c * \frac{v}{\|v\|}.$$

**Camera Look At:** To ensure the object is shown at different locations within the camera frame, the camera Look At needs to be varied. However, range of values such that the object is visible can be present across the entire range of the frame depends on the camera position. So, we sample camera Look At as $l_c$ as follows:

$$l_c \sim \mathrm{Unif}(K * x_c, K * y_c, K * z_c), \text{where } K = 0.3.$$

**Camera Up Vector:** Note that the camera Up Vector is implemented as the vector joining the camera center (0,0,0) to a specified position. We sample this position and therefore the Up Vector $u_c$ as follows:

$$x, y, z \sim \mathrm{Unif}(-1, 1),$$
$$u_c = (x, y, z).$$

**Camera Field of View (FOV):** We sample the field of view $f_c$ while ensuring:

$$f_c \sim \mathrm{Unif}(35, 100).$$

**Light Position:** For every scene we first sample the number of lights $n$ between 1-4 with equal probability. For each light $i$, a random radius $r_i$ is sampled ensuring $r_i \sim \mathrm{Unif}(1, 8)$. Then the light is placed on a random point $(x_i, y_i, z_i)$ on the sphere of radius $r_i$.

**Light Look At:** As above, to ensure that the light is visible on the canvas, light Look At is sampled as a function of the camera position:

$$l_i \sim \mathrm{Unif}(K * x_c, K * y_c, K * z_c), \text{where } K = 0.3.$$

**Light Size:**Every light in our setup is implemented as an area light, and therefore requires a height and width to specify the size. We generate the size $s_i$ for light $i$ as:

$$h, w \sim \text{Unif}(0.1, 5.0),$$
$$s_i = (h, w).$$

**Light Intensity:**This parameter specifies the RGB intensity of the light. For light $i$, RGB color intensity $c_i$ was sampled as:

$$r, g, b \sim \text{Unif}(0, 1),$$
$$c_i = (r, g, b).$$

**Object Material:** To ensure no spurious correlations between object texture and category, all object textures were set to a single diffuse material. More specifically, the material is a linear blend between a Lambertian model and a microfacet model with Phong distribution, with Schilick's Fresnel approximation. Diffuse reflectance was set to 1.0, and the material was set to reflect on both sides.

In Fig. 1 below we show more examples of images from our rendered dataset.

Figure 1: *Sample Image from dataset*

3

## B  Training Details

**Image Normalization:** To get good generalization to unseen 3D models and stable learning, each image was normalized to zero mean and unit standard deviation.

**Batch Size:** All CNN models were trained with a batch size of 75 images, while transformers were trained with a batch size of 25.

**Optimizer and Learning Rate:** All models were trained with an Adam optimizer with a fixed learning rate of 0.0003. Other learning rates including 0.0001, 0.001, 0.01 and 0.1 were tried but they performed either similarly well or worse.

**Number of Epochs:** All models were trained for 50 epochs.

**Hardware:** A compute cluster of NVIDIA TeslaK80 GPUs was used. All models were trained on a single GPU at a time, and a total of 8 GPUs were used across all experiments.

## C  CMA-ES algorithm

We use the CMA-ES algorithm to create our *CMA-Search*, which starts with a correctly classified image and searches the vicinity of the camera or light parameters of the original image for a set of parameters which lead to a misclassification. To implement this, we rely on pycma [2]. The exact subroutines for parameter update are managed automatically by the pycma library. For a detailed overview and theoretical underpinnings of the CMA-ES algorithm, please refer to the tutorial from the authors of CMA-ES [3].

## D  Distribution of camera and light parameters for misclassified image

In Sec.5 of the main text, we showed the camera parameters for misclassified images for one category (Car) and one architecture (ResNet18). As described, errors are distributed across the range of camera and light parameters with no clear patterns in these parameters separating correctly classified and misclassified images. In Fig. 2 below we show that this result is consistent across multiple categories and architectures.

## E  Examples of ImageNet errors discovered by our *CMA-Search* method

In Fig. 3 below, we show more examples of misclassified ImageNet images found using our *CMA-Search* approach. As described in the main text, novel views were generated using the single view MPI model [4], and the results shown here are for the OpenAI CLIP architecture [5].
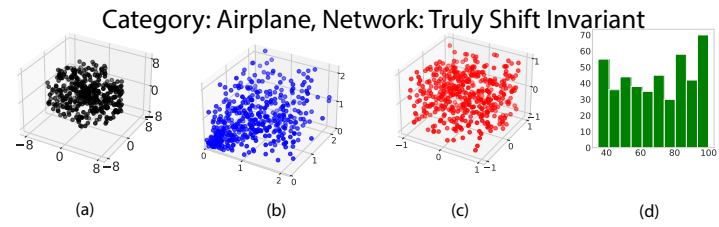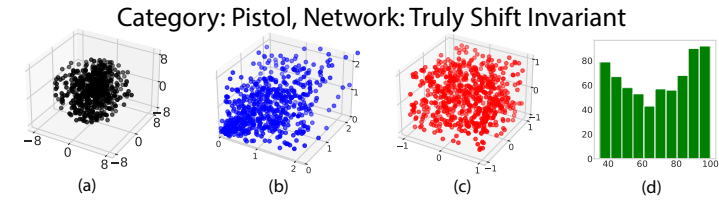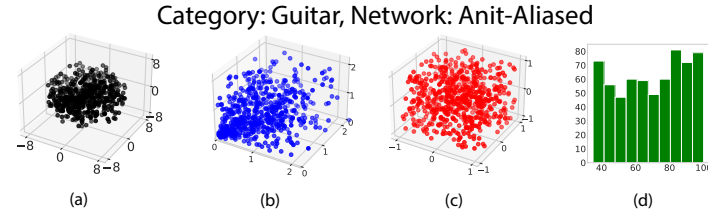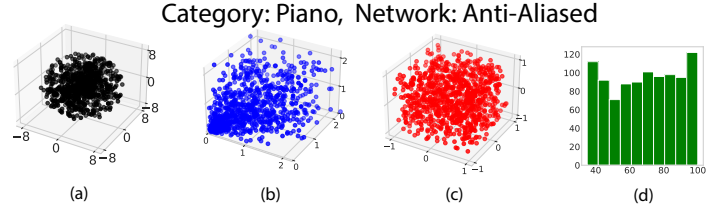
Category: Piano, Network: Anti-Aliased

(a)  (b)  (c)  (d)

Category: Guitar, Network: Anit-Aliased

(a)  (b)  (c)  (d)

Category: Pistol, Network: Truly Shift Invariant

(a)  (b)  (c)  (d)

Category: Airplane, Network: Truly Shift Invariant

(a)  (b)  (c)  (d)

Figure 2: *Camera Parameters that lead to misclassifications for multiple categories and architectures. (a) Camera Position, (b) Camera Look At, (c) Up Vector, (d) Histogram of Lens Field of View.*
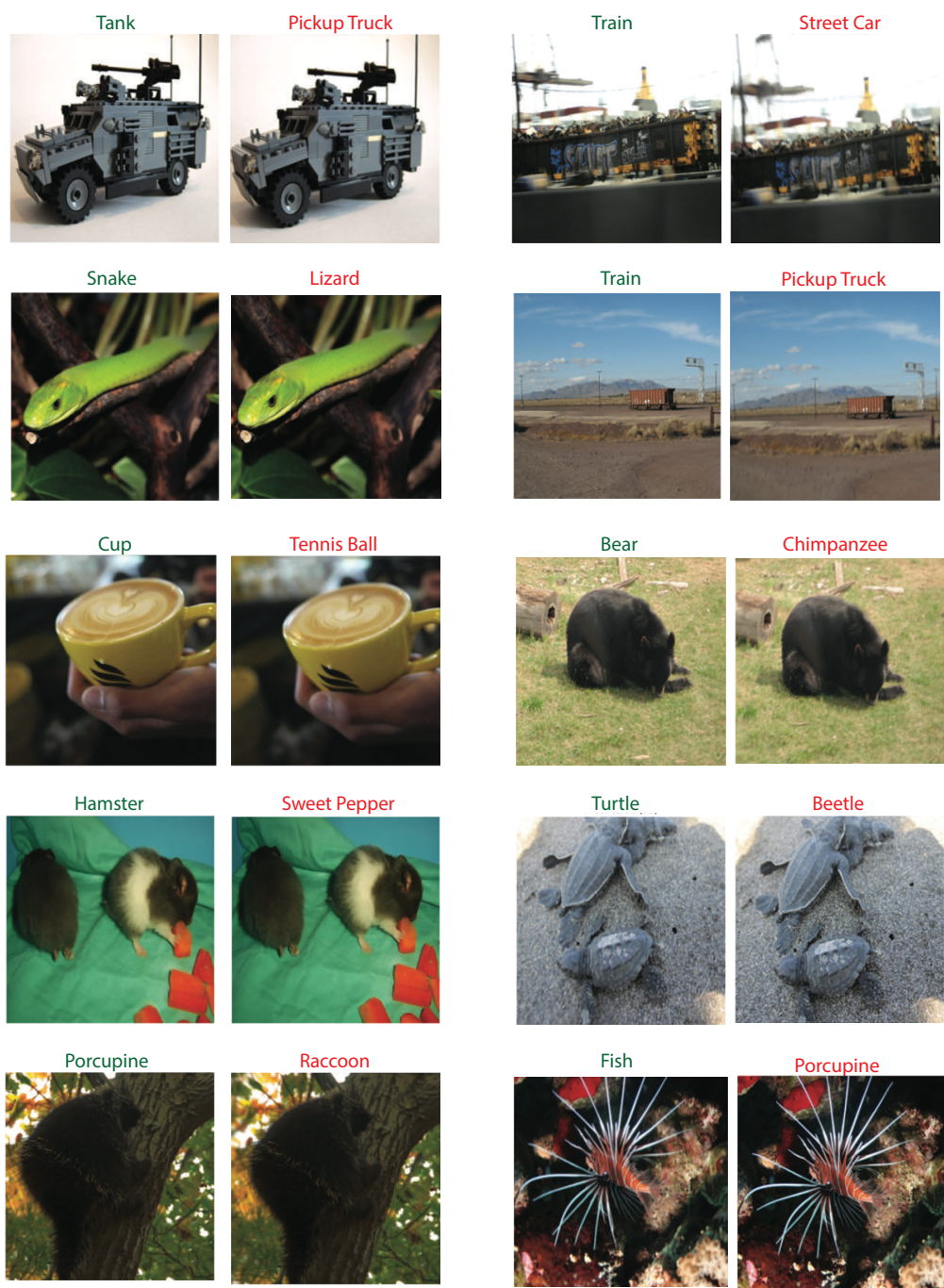
Figure 3: *More examples of misclassified ImageNet images discovered by CMA-Search combined with the single view MPI model.*

# References

[1] Tzu-Mao Li, Miika Aittala, Frédo Durand, and Jaakko Lehtinen. Differentiable Monte Carlo ray tracing through edge sampling. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018.

[2] Nikolaus Hansen, Youhei Akimoto, and Petr Baudis. CMA-ES/pycma on Github. Zenodo, DOI:10.5281/zenodo.2559634, February 2019.

[3] Nikolaus Hansen. The CMA evolution strategy: A tutorial. *arXiv preprint arXiv:1604.00772*, 2016.

[4] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 551–560, 2020.

[5] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.