
Exploiting 3D Shape Bias towards Robust Vision

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Robustness research in machine vision faces a challenge. Many variants of
2 ImageNet-scale robustness benchmarks have been proposed, only to reveal that
3 current vision systems fail under distributional shifts. Although aiming for higher
4 robustness accuracy on these benchmarks is important, we also observe that simply
5 using larger models and larger training datasets may not lead to true robustness,
6 demanding further innovation. To tackle the problem from a new perspective, we
7 encourage closer collaboration between the robustness and 3D vision communities.
8 This proposal is inspired by human vision, which is surprisingly robust to envi-
9 ronmental variation, including both naturally occurring disturbances and artificial
10 corruptions. We hypothesize that such robustness, at least in part, arises from
11 our ability to infer 3D geometry from 2D retinal projections. In this work, we
12 take a first step toward testing this hypothesis by viewing 3D reconstruction as a
13 pretraining method for building more robust vision systems. We introduce a novel
14 dataset called Geon3D, which is derived from objects that emphasize variation
15 across shape features that the human visual system is thought to be particularly
16 sensitive. This dataset enables, for the first time, a controlled setting where we can
17 isolate the effect of “3D shape bias” in robustifying neural networks, and informs
18 new approaches for increasing robustness by exploiting 3D vision tasks. Using
19 Geon3D, we find that CNNs pretrained on 3D reconstruction are more resilient to
20 viewpoint change, rotation, and shift than regular CNNs. Further, when combined
21 with adversarial training, 3D reconstruction pretrained models improve adversarial
22 and common corruption robustness over vanilla adversarially-trained models. We
23 hope that our findings and dataset will encourage exploitation of synergies between
24 the robustness researchers, 3D computer vision community, and computational
25 perception researchers in cognitive science, paving a way for achieving human-like
26 robustness under complex, real-world stimuli conditions.

27 1 Introduction

28 Building robust vision systems is a major open problem. Tremendous efforts have been made since
29 adversarial examples were first reported [36], and yet adversarial robustness remains perhaps the most
30 important challenge in safe, real-world deployment of modern computer vision systems. Ensuring
31 robustness against more common distributional shifts such as blur and snow also remains a significant
32 challenge [18]. As clean ImageNet accuracy saturates, the research community has developed various
33 ImageNet-scale benchmarks to evaluate the performance of vision models under distributional shifts
34 such as broader viewpoint variability [3], style and texture change [15], geographic shifts [19].
35 These benchmarks, as well as the recent algorithms that are evaluated using smaller-scale datasets
36 such as MNIST and CIFAR10 [38, 39], reveal that current vision systems have plenty of room for
37 improvement in terms of robustness.

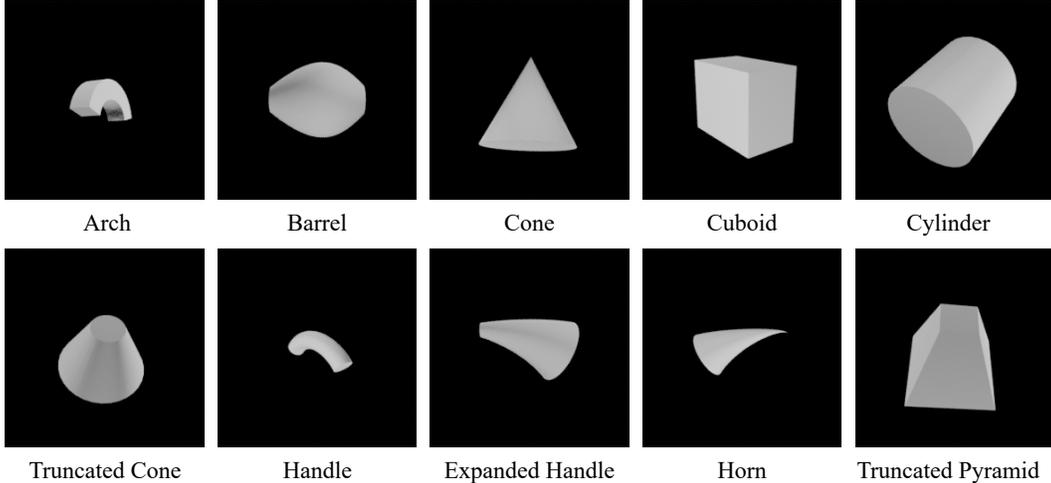


Figure 1: Examples of 10 Geon categories from Geon3D-10. The full list of 40 Geons we construct (Geon3D-40) is provided in the Appendix.

38 So far, robustness research in machine vision focuses on classification. Models trained for image
 39 classification might learn to associate class labels with a limited range of surface-related cues such
 40 as image contours, but they do not fully or explicitly reflect the relationship between 3D objects
 41 and how they are projected to images. On the contrary, the human visual system recovers rich
 42 three-dimensional (3D) geometry, including objects, shapes and surfaces, from two-dimensional
 43 (2D) retinal inputs. This ability to make inferences about the underlying scene structure from
 44 input images—also known as analysis-by-synthesis—is thought to be critical for the robustness of
 45 biological vision to occlusions, distortions, and lighting variations [41, 26].

46 While aiming for higher accuracy on ImageNet-scale benchmarks is important, the current landscape
 47 of robustness research shows that we face a clear challenge [37]. In fact, the consensus seems to be that
 48 large models and large training data work well for some distribution shifts, but nothing consistently
 49 help in all variants of ImageNet robustness benchmarks, awaiting methodological innovation to
 50 achieve human-level robustness [19]. To unblock the situation, we advocate closer collaboration
 51 between the robustness and 3D vision communities, in the hope of fostering new types of robustness
 52 research. This paper serves as a first step towards this effort, where we focus on learning features
 53 to facilitate inferences about 3D object shape. Our goal is to test the hypothesis that shape bias—
 54 learning representations that enable accurate inferences of 3D from 2D, which we refer to as “3D
 55 shape bias”—will induce robustness to naturally occurring challenging viewing conditions (e.g., fog,
 56 snow, brightness) and artificial image corruptions (e.g., due to adversarial attacks).

57 To achieve this, we introduce *Geon3D*—a novel dataset comprised of simple yet realistic shape
 58 variations, derived from the human object recognition hypothesis called Geon Theory [5]. This
 59 dataset enables us to study, in a controlled setting, 3D shape bias of 3D reconstruction models
 60 that learn to represent shapes solely from 2D supervision [28]. We find that CNNs trained for 3D
 61 reconstruction are more robust to unseen viewpoints, rotation and translation than regular CNNs.
 62 Moreover, when combined with adversarial training, 3D reconstruction pretraining improves common
 63 corruption and adversarial robustness over CNNs that only use adversarial training. These results
 64 suggest that the Geon3D dataset provides a controlled and effective measure of robustness, and unlike
 65 existing, commonly used datasets in this area such as CIFAR10 and ImageNet-C, Geon3D guides
 66 novel approaches by facilitating an interface between robust machine learning and 3D reconstruction.
 67 (Please see the Related Work section for a discussion of Geon3D in the context of existing 3D shape
 68 datasets.)

69 Biological vision is not only about object classification or localization, but also about making rich
 70 inference about the underlying causes of scenes such as 3D shapes and surfaces [29, 41, 26]. We hope
 71 our findings and dataset will encourage the community to tackle robustness problems through the
 72 lens of 3D inference and the perspective of perception as analysis-by-synthesis, toward the combined
 73 goals of building machine vision systems with human-like richness and reliability.

74 2 Approach

75 We first describe the Geon Theory, which our dataset construction relies on. Next, we explain the
76 data generation process used in the creation of Geon3D (§2.1), and how we train a 3D reconstruction
77 model (§2.2).

78 2.1 Geon3D Benchmark

79 The concept of *Geons*—or *Geometric ions*—was originally introduced by Biederman as the building
80 block for his Recognition-by-Components (RBC) Theory [5]. The RBC theory argues that human
81 shape perception segments an object at regions of sharp concavity, modeling an object as a com-
82 position of Geons—a subset of generalized cylinders [6]. Similar to generalized cylinders, each
83 Geon is defined by its axis function, cross-section shape, and sweep function. In order to reduce
84 the possible set of generalized cylinders, Biederman considered the properties of the human visual
85 system. He noted that the human visual system is better at distinguishing between straight and curved
86 lines than at estimating curvature; detecting parallelism than estimating the angle between lines; and
87 distinguishing between vertex types such as an arrow, Y, and L-junction [21].

Table 1: Latent features of Geons. S: Straight, C: Curved, Co: Constant, M: Monotonic, EC: Expand and Contract, CE: Contract and Expand, T: Truncated, P: End in a point, CS: End as a curved surface

Feature	Values
Axis	S, C
Cross-section	S, C
Sweep function	Co, M, EC, CE
Termination	T, P, CS

Table 2: Similar Geon categories, where only a single feature differs out of four shape features. “T.” stands for “Truncated”. “E.” stands for “Expanded”.

Geon Category	Difference
Cone vs. Horn	Axis
Handle vs. Arch	Cross-section
Cuboid vs. Cylinder	Cross-section
T. Pyramid vs. T. Cone	Cross-section
Cuboid vs. Pyramid	Sweep function
Barrel vs. T. Cone	Sweep function
Horn vs. E. Handle	Termination

88 Our focus in this paper is not the RBC theory or whether it is the right way to think about how we see
89 shapes. Instead, we wish to build upon the way Biederman characterized these Geons. Biederman
90 proposed using two to four values to characterize each feature of Geons. Namely, the axis can be
91 straight or curved; the shape of cross section can be straight-edged or curved-edged; the sweep
92 function can be constant, monotonically increasing / decreasing, monotonically increasing and then
93 decreasing (i.e. expand and contract), or monotonically decreasing and then increasing (i.e. contract
94 and expand); the termination can be truncated, end in a point, or end as a curved surface. A summary
95 of these dimensions is given in Table 1.

96 Representative Geon classes are shown in Figure 1. For example, the “Arch” class is uniquely
97 characterized by its curved axis, straight-edged cross section, constant sweep function, and truncated
98 termination. These values of Geon features are *nonaccidental*—we can determine whether the axis is
99 straight or curved from almost any viewpoint, except for a few *accidental* cases. For instance, an
100 arch-like curve in the 3D space is perceived as a straight line only when the viewpoint is aligned in a
101 way that the curvature vanishes. These properties make Geons an ideal dataset to analyze 3D shape
102 bias and part-level robustness of vision models. For details of data preparation, see Appendix.

103 2.2 3D reconstruction as pretraining

104 To explore advantages of direct approaches to induce shape bias in vision models, we turn our
105 attention to a class of 3D reconstruction models. The main hypothesis of our study is that the task of
106 3D reconstruction pressures the model to obtain robust representations.

107 Recently, there has been significant progress in learning-based approaches to 3D reconstruction,
108 where the data representation can be classified into voxels [10, 32], point clouds [14, 1], mesh [22, 17],
109 and neural implicit representations [25, 9, 31, 35]. We focus on neural implicit representations, where
110 models learn to implicitly represent 3D geometry in neural network parameters after training. We
111 avoid models that require 3D supervision such as ground truth 3D shapes. This is because we are

112 interested in models that only require 2D supervision for training and how inductive bias of 2D-to-3D
113 inference achieves robustness.

114 Specifically, we use Differentiable Volumetric Rendering (DVR) [28], which consists of a CNN-based
115 image encoder and a differentiable neural rendering module. We train DVR to reconstruct 3D shapes
116 of Geon3D-10. For more details of DVR and 3D reconstruction, we refer the readers to the original
117 paper [28].

118 3 Experimental Results

119 In this section, we demonstrate how 3D shape bias improves model robustness on the Geon3D-10
120 classification under various image perturbations. Our 3D-shape-biased classifier is based on the image
121 encoder of the 3D reconstruction model (DVR) that is pretrained to reconstruct Geon3D-10. We add
122 a linear classification layer on top of the image encoder, and then finetune, either just that linear layer
123 (**DVR-Last**) or the entire encoder (**DVR**), for Geon3D-10 classification. Our baseline is a vanilla
124 neural network (**Regular**) that is trained normally for Geon3D-10 classification. To see the difference
125 between 3D shape bias and 2D shape bias in the sense of [15], we also evaluate the following models,
126 which are hypothesized to rely their prediction more on shape than texture. **Stylized** is a model
127 trained on Stylized images of Geons. **Adversarially trained network (AT)** is a network that uses
128 adversarial examples during training [24]. **InfoDrop** [34] is a recently proposed model that induces
129 2D shape bias by decorrelating each layer’s output with texture. To control for variation in network
130 architectures, we use ImageNet-pretrained ResNet18 for all models we tested. The image encoder of
131 DVR is also initialized using ImageNet-pretrained training for 3D reconstruction of Geons.

132 **Background variations** To quantify the effect of textured background, we prepare three versions
133 of Geon3D-10: black background, random textured background (Geon3D-10-RandTextured), and
134 correlated background (Geon3D-10-CorrTextured). For Geon3D-10-RandTextured, we replace
135 each black background with a random texture image out of 10 texture categories chosen from the
136 Describable Textures Dataset (DTD) [11]. For Geon3D-10-CorrTextured, we choose 10 texture
137 categories from DTD and introduce spurious correlations between Geon category and texture class
138 (i.e., each Geon category is paired with one texture class). Examples of Geon3D with textured
139 background are shown in Figure 4 (Right). These three versions of our dataset allow us to analyze
140 more realistic image conditions as well as to test robustness despite variation and distributional shifts
141 in textures.

142 **Accuracy under rotation and translation (shifting pixels)** CNNs are known to be vulnerable to
143 rotation and shifting of the image pixels [2]. As shown in Table 3, our model (DVR) pretrained with
144 3D reconstruction performs better than all other models under rotation and shift even though it is not
145 explicitly trained to defend against those attacks. We observe that DVR-Last performs second best,
146 indicating that this “for free” robustness to rotation and shift is largely in place even when finetuning
147 on the classification task is restricted to only linear decoding of the categories.

Table 3: Accuracy of shape-biased classifiers against rotation and shifting of pixels on Geon3D under
unseen viewpoints. We randomly add rotations of at most 30° and translations of at most 10% of the
image size in each x, y direction. We report the mean accuracy and standard deviation over 5 runs of
this stochastic procedure over the entire evaluation set.

	REGULAR	INFODROP	STYLIZED	AT- L_2	AT- L_∞	DVR-LAST	DVR
ROTATION	82.18 _(1.06)	80.76 _(0.69)	78.47 _(0.57)	87.00 _(0.57)	89.58 _(0.48)	90.44 _(0.30)	93.46 _(0.44)
SHIFT	72.28 _(0.43)	71.86 _(0.63)	61.44 _(0.29)	53.84 _(0.71)	61.50 _(1.11)	73.24 _(0.73)	76.52 _(0.89)

148 3.1 Robustness against Common Corruptions

149 In this section, we show that, when combined with adversarial training, 3D pretrained models
150 (denoted as DVR+AT- L_2 and DVR+AT- L_∞) improve robustness against common image corruptions,
151 above and beyond what can be accomplished just using adversarial training. For these models, we
152 use adversarial training during the finetuning of the 3D reconstruction model for the Geon3D-10
153 classification task. Here we evaluate the effect of 3D shape bias not only in the somewhat sterile

154 scenario of the clean, black background images, but also using the background-textured versions
 155 of our dataset. To do this, we train all models using Geon3D-10-RandTextured, where we replace
 156 the black background with textures randomly sampled from DTD (see Figure 4, right panel, for
 157 examples). During evaluation, we use unseen viewpoints.

158 The results are shown in Table 4. We see that starting adversarial training from DVR-pretrained
 159 weights improves robustness across all corruption types, over what can be achieved by only either
 160 $AT-L_2$ or $AT-L_\infty$. DVR-AT and AT models fail on ‘‘Contrast’’ and ‘‘Fog’’. This has been a known
 161 issue for AT [16], which requires future work to explore. While Stylized performs best under certain
 162 corruption types, we can see that DVR- $AT-L_2$ leads to broader robustness across the corruptions we
 163 considered.

Table 4: Accuracy of classifiers against common corruptions under unseen viewpoints. All models are trained and evaluated on Geon3D-10 with random textured background. Pretraining on 3D shape reconstruction using DVR leads to broader robustness relative to other models.

	REGULAR	INFODROP	STYLIZED	$AT-L_2$	$AT-L_\infty$	DVR+ $AT-L_2$	DVR+ $AT-L_\infty$
INTACT	0.741	0.596	0.701	0.691	0.464	0.758	0.513
PIXELATE	0.608	0.458	0.653	0.623	0.415	0.719	0.470
DEFOCUS BLUR	0.154	0.152	0.402	0.490	0.298	0.605	0.349
GAUSSIAN NOISE	0.222	0.465	0.601	0.555	0.412	0.701	0.470
IMPULSE NOISE	0.187	0.270	0.497	0.322	0.136	0.594	0.148
FROST	0.144	0.269	0.638	0.142	0.209	0.148	0.240
FOG	0.338	0.281	0.659	0.187	0.120	0.264	0.130
ELASTIC	0.427	0.314	0.428	0.416	0.266	0.499	0.307
JPEG	0.414	0.422	0.634	0.629	0.434	0.731	0.484
CONTRAST	0.408	0.286	0.673	0.141	0.120	0.179	0.135
BRIGHTNESS	0.525	0.518	0.702	0.500	0.388	0.549	0.429
ZOOM BLUR	0.334	0.238	0.560	0.518	0.327	0.639	0.378

164 3.2 3D Pretraining Improves Adversarial Robustness

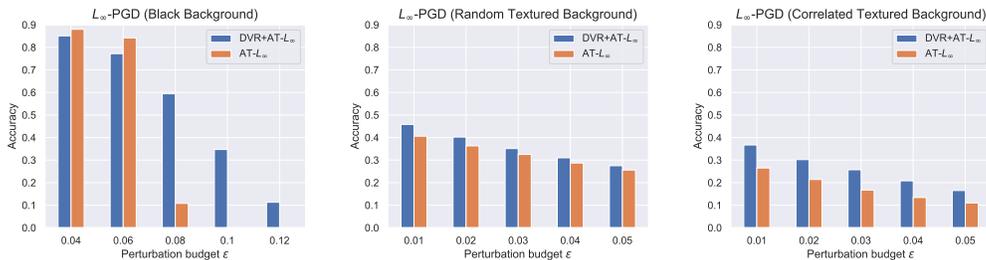


Figure 2: Robustness comparison between $AT-L_\infty$ and $DVR+AT-L_\infty$ with increasing perturbation budget ϵ on three variations of Geon3D-10. We use L_∞ -PGD with 100 iterations and $\epsilon/10$ to be the stepsize. See Appendix for $AT-L_2$ results, where we also find that 3D pretraining improves vanilla AT models.

165 In this section, we show that 3D pretrained AT models improve adversarial robustness over vanilla AT
 166 models. We attack our models using L_∞ -PGD [24], with 100 iterations and $\epsilon/10$ to be the stepsize,
 167 where ϵ is the perturbation budget. We compare $AT-L_\infty$ and $DVR+AT-L_\infty$ for black, randomly
 168 textured, and correlated textured backgrounds. The results are shown in Figure 2. In the black
 169 background set, while 3D pretrained AT slightly performs worse than vanilla AT for smaller epsilon
 170 values, it significantly robustifies AT-trained models for large epsilons. A small but appreciable gain
 171 in robustness can be seen for the other two backgrounds types. These pattern of results are consistent
 172 across attack types, with DVR providing significant robustness over vanilla AT under the L_2 regime
 173 (see Appendix).

174 References

175 [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning Representations and Generative Models for 3D Point Clouds. In *International Conference on Machine Learning*, pp. 40–49. PMLR, July 2018.
 176
 177

- 178 [2] Aharon Azulay and Yair Weiss. Why do deep convolutional networks generalize so poorly to
179 small image transformations? *Journal of Machine Learning Research*, pp. 25, 2019.
- 180 [3] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund,
181 Josh Tenenbaum, and Boris Katz. ObjectNet: A large-scale bias-controlled dataset for pushing
182 the limits of object recognition models. *Advances in Neural Information Processing Systems*,
183 32:9453–9463, 2019.
- 184 [4] Barr. Superquadrics and Angle-Preserving Transformations. *IEEE Computer Graphics and
185 Applications*, 1(1):11–23, January 1981. ISSN 1558-1756. doi: 10.1109/MCG.1981.1673799.
- 186 [5] Irving Biederman. Recognition-by-components: A theory of human image understanding.
187 *Psychological Review*, 94(2):115–147, 1987. ISSN 1939-1471(Electronic),0033-295X(Print).
188 doi: 10.1037/0033-295X.94.2.115.
- 189 [6] I. Binford. Visual Perception by Computer. *IEEE Conference of Systems and Control*, 1971.
- 190 [7] Online Community Blender. Blender - a 3D modelling and rendering package. Blender
191 Foundation, 2021.
- 192 [8] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo
193 Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher
194 Yu. ShapeNet: An Information-Rich 3D Model Repository. *arXiv:1512.03012 [cs]*, December
195 2015.
- 196 [9] Zhiqin Chen and Hao Zhang. Learning Implicit Fields for Generative Shape Modeling. In *2019
197 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5932–5941,
198 Long Beach, CA, USA, June 2019. IEEE. ISBN 978-1-72813-293-8. doi: 10.1109/CVPR.2019.
199 00609.
- 200 [10] C. Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and S. Savarese. 3D-R2N2: A Unified
201 Approach for Single and Multi-view 3D Object Reconstruction. In *ECCV*, 2016. doi: 10.1007/
202 978-3-319-46484-8_38.
- 203 [11] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi.
204 Describing Textures in the Wild. In *2014 IEEE Conference on Computer Vision and Pattern
205 Recognition*, pp. 3606–3613, Columbus, OH, USA, June 2014. IEEE. ISBN 978-1-4799-5118-5.
206 doi: 10.1109/CVPR.2014.461.
- 207 [12] P. Dayan, Geoffrey E. Hinton, R. Neal, and R. Zemel. The Helmholtz Machine. *Neural
208 Computation*, 1995. doi: 10.1162/neco.1995.7.5.889.
- 209 [13] S. M. Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S. Morcos, Marta
210 Garnelo, Avraham Ruderman, Andrei A. Rusu, Ivo Danihelka, Karol Gregor, David P. Reichert,
211 Lars Buesing, Theophane Weber, Oriol Vinyals, Dan Rosenbaum, Neil Rabinowitz, Helen
212 King, Chloe Hillier, Matt Botvinick, Daan Wierstra, Koray Kavukcuoglu, and Demis Hassabis.
213 Neural scene representation and rendering. *Science*, 360(6394):1204–1210, June 2018. ISSN
214 0036-8075, 1095-9203. doi: 10.1126/science.aar6170.
- 215 [14] Haoqiang Fan, Hao Su, and Leonidas J. Guibas. A Point Set Generation Network for 3D Object
216 Reconstruction From a Single Image. In *Proceedings of the IEEE Conference on Computer
217 Vision and Pattern Recognition*, pp. 605–613, 2017.
- 218 [15] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and
219 Wieland Brendel. ImageNet-trained CNNs are biased towards texture; increasing shape bias
220 improves accuracy and robustness. In *International Conference on Learning Representations*,
221 September 2018.
- 222 [16] Justin Gilmer, Nicolas Ford, Nicholas Carlini, and Ekin Cubuk. Adversarial Examples Are a
223 Natural Consequence of Test Error in Noise. In *International Conference on Machine Learning*,
224 pp. 2280–2289. PMLR, May 2019.

- 225 [17] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, and Mathieu Aubry.
226 A Papier-Mâché Approach to Learning 3D Surface Generation. In *Proceedings of the IEEE*
227 *Conference on Computer Vision and Pattern Recognition*, pp. 216–224, 2018.
- 228 [18] Dan Hendrycks and Thomas Dietterich. Benchmarking Neural Network Robustness to Common
229 Corruptions and Perturbations. In *International Conference on Learning Representations*,
230 September 2018.
- 231 [19] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,
232 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and
233 Justin Gilmer. The Many Faces of Robustness: A Critical Analysis of Out-of-Distribution
234 Generalization. *ICCV*, 2021.
- 235 [20] X. Huang and S. Belongie. Arbitrary Style Transfer in Real-Time with Adaptive Instance
236 Normalization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 1510–
237 1519, October 2017. doi: 10.1109/ICCV.2017.167.
- 238 [21] Katsushi Ikeuchi (ed.). *Computer Vision: A Reference Guide*. Springer US, 2014. ISBN
239 978-0-387-30771-8.
- 240 [22] H. Kato, Y. Ushiku, and T. Harada. Neural 3D Mesh Renderer. In *2018 IEEE/CVF Conference*
241 *on Computer Vision and Pattern Recognition*, pp. 3907–3916, June 2018. doi: 10.1109/CVPR.
242 2018.00411.
- 243 [23] T. D. Kulkarni, P. Kohli, J. B. Tenenbaum, and V. Mansinghka. Picture: A probabilistic
244 programming language for scene perception. In *2015 IEEE Conference on Computer Vision and*
245 *Pattern Recognition (CVPR)*, pp. 4390–4399, June 2015. doi: 10.1109/CVPR.2015.7299068.
- 246 [24] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
247 Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference*
248 *on Learning Representations*, February 2018.
- 249 [25] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger.
250 Occupancy Networks: Learning 3D Reconstruction in Function Space. In *2019 IEEE/CVF*
251 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4455–4465, Long Beach,
252 CA, USA, June 2019. IEEE. ISBN 978-1-72813-293-8. doi: 10.1109/CVPR.2019.00459.
- 253 [26] David Mumford. Pattern Theory: A Unifying Perspective. In Anthony Joseph, Ful-
254 bert Mignot, François Murat, Bernard Prum, and Rudolf Rentschler (eds.), *First European*
255 *Congress of Mathematics: Paris, July 6-10, 1992 Volume I Invited Lectures (Part 1)*, Progress
256 in Mathematics, pp. 187–224. Birkhäuser, Basel, 1994. ISBN 978-3-0348-9110-3. doi:
257 10.1007/978-3-0348-9110-3_6.
- 258 [27] Chaithanya Kumar Mummadi, Ranjitha Subramaniam, Robin Huttmacher, Julien Vitay, Volker
259 Fischer, and Jan Hendrik Metzen. Does enhanced shape bias improve neural network robustness
260 to common corruptions? In *International Conference on Learning Representations*, September
261 2020.
- 262 [28] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable
263 Volumetric Rendering: Learning Implicit 3D Representations Without 3D Supervision. In *2020*
264 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3501–3512,
265 Seattle, WA, USA, June 2020. IEEE. ISBN 978-1-72817-168-5. doi: 10.1109/CVPR42600.
266 2020.00356.
- 267 [29] Bruno A Olshausen. Perception as an Inference Problem. *The Cognitive Neurosciences, Sixth*
268 *Edition | The MIT Press*, pp. 18, 2013.
- 269 [30] Stephen E. Palmer. *Vision Science : Photons to Phenomenology*. MIT Press, 1999.
- 270 [31] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove.
271 DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation. In *2019*
272 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 165–174,
273 Long Beach, CA, USA, June 2019. IEEE. ISBN 978-1-72813-293-8. doi: 10.1109/CVPR.2019.
274 00025.

- 275 [32] G. Riegler, A. O. Ulusoy, and A. Geiger. OctNet: Learning Deep 3D Representations at High
276 Resolutions. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
277 pp. 6620–6629, July 2017. doi: 10.1109/CVPR.2017.701.
- 278 [33] Lukas Schott, Jonas Rauber, Matthias Bethge, and Wieland Brendel. Towards the first adver-
279 sarily robust neural network model on MNIST. In *International Conference on Learning*
280 *Representations*, September 2018.
- 281 [34] Baifeng Shi, Dinghui Zhang, Qi Dai, Zhanxing Zhu, Yadong Mu, and Jingdong Wang. Infor-
282 mative Dropout for Robust Representation Learning: A Shape-bias Perspective. In *International*
283 *Conference on Machine Learning*, pp. 8828–8839. PMLR, November 2020.
- 284 [35] Vincent Sitzmann, Michael Zollhoefer, and Gordon Wetzstein. Scene Representation Net-
285 works: Continuous 3D-Structure-Aware Neural Scene Representations. In *Advances in Neural*
286 *Information Processing Systems*, pp. 1121–1132, 2019.
- 287 [36] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Good-
288 fellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference*
289 *on Learning Representations*, 2014.
- 290 [37] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, B. Recht, and Ludwig Schmidt.
291 Measuring Robustness to Natural Distribution Shifts in Image Classification. *NeurIPS*, 2020.
- 292 [38] Florian Tramer and Dan Boneh. Adversarial Training and Robustness for Multiple Perturbations.
293 In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.,
294 2019.
- 295 [39] Hanshu Yan, Jingfeng Zhang, Gang Niu, Jiashi Feng, Vincent Tan, and Masashi Sugiyama.
296 CIFS: Improving Adversarial Robustness of CNNs via Channel-wise Importance-based Feature
297 Selection. In *Proceedings of the 38th International Conference on Machine Learning*, pp.
298 11693–11703. PMLR, July 2021.
- 299 [40] Ilker Yildirim, Mario Belledonne, Winrich Freiwald, and Josh Tenenbaum. Efficient inverse
300 graphics in biological face processing. *Science Advances*, 6(10):eaax5979, March 2020. ISSN
301 2375-2548. doi: 10.1126/sciadv.aax5979.
- 302 [41] Alan Yuille and Daniel Kersten. Vision as Bayesian inference: Analysis by synthesis? *Trends*
303 *in Cognitive Sciences*, 10(7):301–308, July 2006. ISSN 1364-6613. doi: 10.1016/j.tics.2006.05.
304 002.
- 305 [42] Tianyuan Zhang and Zhanxing Zhu. Interpreting Adversarially Trained Convolutional Neural
306 Networks. In *International Conference on Machine Learning*, pp. 7502–7511. PMLR, May
307 2019.
- 308 [43] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and
309 Jianxiong Xiao. 3D ShapeNets: A deep representation for volumetric shapes. In *2015 IEEE*
310 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1912–1920, Boston, MA,
311 USA, June 2015. IEEE. ISBN 978-1-4673-6964-0. doi: 10.1109/CVPR.2015.7298801.

312 A Additional experiments

313 A.1 3D shape bias improves generalization to unseen views and reduces similar category 314 confusion

315 One of the crucial but often overlooked examples of 3D shape bias that human vision has is “visual
316 completion” [30], which refers to our ability to infer portions of surface that we cannot actually see.
317 For instance, when we look at the top-left image in Figure 4, we automatically recognize it as a whole
318 cube, even though we cannot see its rear side. We view the task of 3D reconstruction as a way to
319 build such an ability into neural networks. In this section, we investigate how such 3D shape bias of
320 DVR improves classification of similar Geon categories under unseen viewpoints, testing both DVR
321 (where we finetune all layers of the image encoder) and DVR-Last (where we finetune only the top
322 classification layer of the image encoder).

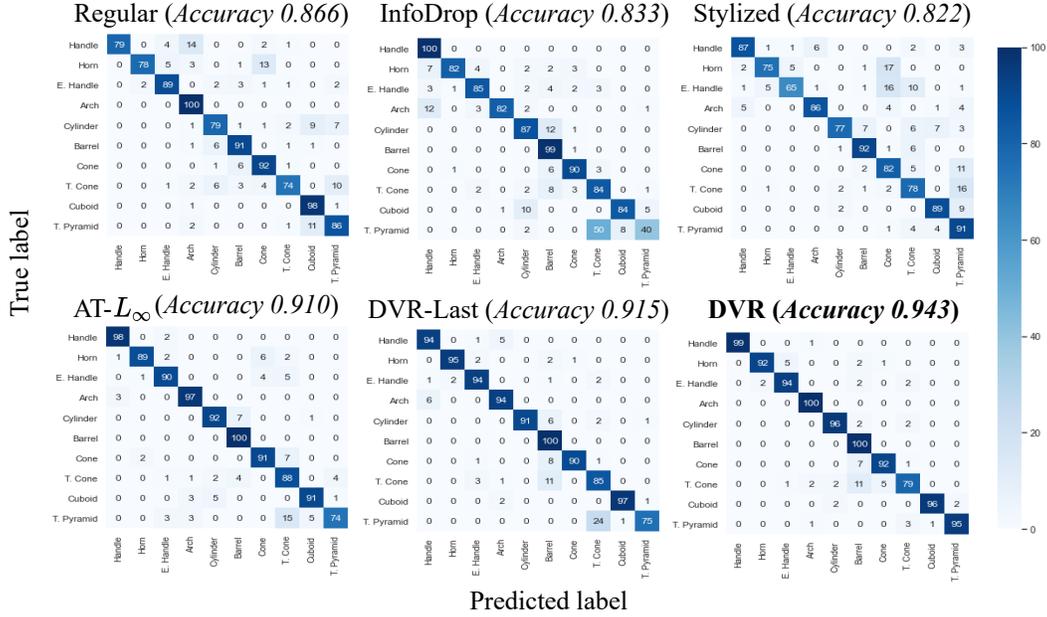


Figure 3: Accuracy per Geon category under unseen viewpoints. Even though all models perform reasonably well, there is still a range of overall accuracy values. In addition, we see that when networks make a mistake, it is often between similar Geon categories (see Table 2 for a list of similar Geon categories). Regular: a baseline model; InfoDrop: a shape-biased model; AT: adversarially trained; Stylized: a network trained on “stylized” version of Geon3D; DVR: We use pretrained weights of the image encoder of Differentiable Volumetric Rendering (3D reconstruction model), a 3D reconstruction model, and finetune all of its layers on the Geon3D-10 classification task. DVR-Last refers to the version where we finetune only the last classification layer.

323 The results of per-category classification are shown in Figure 3. We say two Geons are similar when
 324 there is only a single shape feature difference, as summarized in Table 2. We see that networks often
 325 misclassify similar Geon categories. The vanilla neural network (Regular) often misclassifies “Cone”
 326 vs. “Horn”, “Handle” vs. “Arch”, “Cuboid” vs. “Truncated pyramid”, as well as “Truncated cone” vs.
 327 “Truncated pyramid”. The Geon pairs the InfoDrop model misclassifies include: “Arch” vs. “Handle”,
 328 “Cylinder” vs. “Barrel”, “Cuboid” vs. “Cylinder” and “Truncated pyramid” vs. “Truncated cone”,
 329 which are all pairs with single shape feature difference.

330 Notably, the Stylized model, which is hypothesized to increase bias towards shape-related features,
 331 makes a number of mistakes for similar Geon classes (i.e. “Horn” vs. “Cone”, “Cone” vs. “Truncated
 332 pyramid”, and “Truncated cone” vs. “Truncated pyramid”), similar to the Regular model. This result
 333 is consistent with the finding that the Stylized approach [15] does not necessarily induce proper shape
 334 bias [27].

335 AT- L_∞ and DVR-Last perform better than the models listed above, yet still struggle to distinguish
 336 “Truncated Pyramid” from “Truncated Cone”, where the difference is whether the cross-section
 337 is curved or straight (see Table 2). On the other hand, DVR successfully distinguishes these two
 338 categories. This shows that 3D pretraining before finetuning for the task of classification facilitates
 339 recognition of even highly similar shapes. The hardest pair for DVR is “Truncated cone” vs. “Barrel”,
 340 but the errors the model make appear sensible (Figure 4, middle panel): For example, when the camera
 341 points at the smaller side of the “Truncated Cone”, then there is uncertainty whether the surface
 342 extends beyond self-occlusion by contracting (which would be consistent with the “Barrel” category)
 343 or the surface ends at the point of self-occlusion (which would be consistent with the category
 344 “Truncated Cone”). Indeed, when we inspected the samples of “Truncated Cone” misclassified as
 345 “Barrel” by DVR, we found that for half of those images, the larger side of “Truncated Cone” was
 346 self-occluded. Future psychophysical work should quantitatively compare errors made by these
 347 models to human behavior.

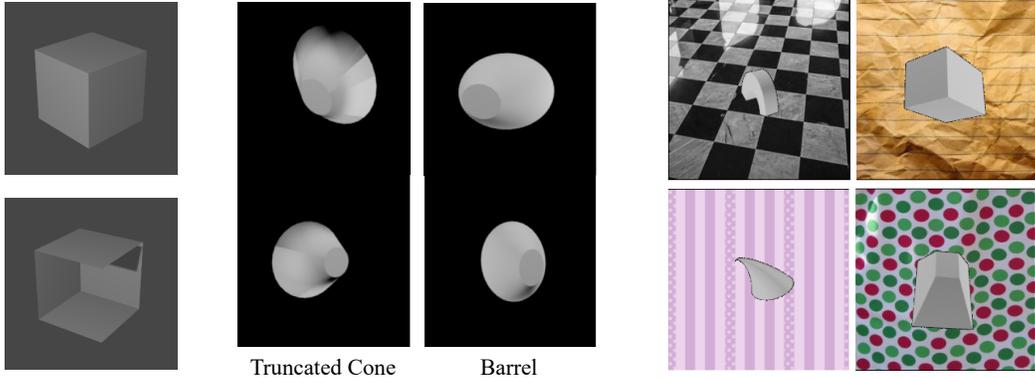


Figure 4: (Left) We humans recognize the top image as a whole cube, automatically filling in the surfaces of its rear, invisible side, although, in principle, there are infinitely many scenes consistent with the sense data, one of which is shown in the bottom image [30]. This illustrates that certain shapes are more readily perceived by the human visual system than others. (Middle) Examples of “Truncated Cone” that are misclassified as “Barrel” by DVR, next to “Barrel” exemplars shown at similar viewpoints. (Right) Example images from Geon3D-10 with textured backgrounds.

348 A.2 Robustness to Distributional Shift in Backgrounds

349 In this section, we evaluate network’s robustness to distributional shift in backgrounds. To do
 350 this, we train all the models on Geon3D-10-CorrTextured, where we introduce spurious correlation
 351 between textured background and Geon category. Therefore, during training, a model can pick up
 352 classification signal from both the shape of Geon as well as background texture. To evaluate trained
 353 models for background shift, we prepare a test set that breaks the correlation between Geon category
 354 and background texture class by cyclically shifting the texture class from i to $i + 1$ for $i = 0, \dots, 9$,
 355 where the class 10 is mapped to the class 0. This is inspired by [15], where they create shape-texture
 356 conflicts to measure 2D shape bias in networks trained for ImageNet classification. However, in our
 357 case, distributional shift from training to test set is designed to isolate and better measure shape bias
 358 by fully disentangling the contributions of texture and shape.

359 The results are shown in Table 5. We see that 2D shape biased models all perform worse than the
 360 3D shape-biased model (DVR+AT- L_∞). Combining AT with 3D pretraining improves classification
 361 accuracy more than 10 % with respect to the best performing variant of AT.

362 Interestingly, comparing randomized vs. correlated background experiments reveals a stark difference
 363 between the two commonly used perturbations in adversarial training (L_2 vs. L_∞). Unlike our
 364 analysis with uncorrelated, randomized backgrounds, we find that adversarial training using L_2 norm
 365 completely biases the model towards texture (no apparent shape bias) when such spurious correlation
 366 between texture and shape category exists.

Table 5: Accuracy of shape-biased classifiers against distributional shift in backgrounds. Here, all models are trained on Geon3D-10-CorrTextured (with background textures correlated with shape categories) and evaluated on a test set where we break this correlation. See Appendix for results using other common corruptions, where we find DVR+AT- L_∞ provides broadest robustness across the corruptions we tested.

REGULAR	INFODROP	STYLIZED	AT- L_2	AT- L_∞	DVR+AT- L_2	DVR+AT- L_∞
0.045	0.121	0.268	0.015	0.311	0.219	0.439

367 B How important is 3D inference?

368 In this section, we investigate the importance of causal 3D inference to obtain good representations.
 369 That is, we explore the impact of having an actual rendering function constrain the representations
 370 learned by a model. Our goal in this section is not to further evaluate the robustness of these features,

371 but to measure the efficiency of representations learned under the constraint of a rendering function
 372 for the basic task of classification.

373 To isolate this effect, we compare DVR to Generative Query Networks (GQN) [13]—a scene
 374 representation model that can generate scenes from unobserved viewpoints—on novel exemplars
 375 from the Geon3D-10 dataset, but using views seen during training. The crucial difference between
 376 DVR and GQN is that GQN does not model the geometry of the object explicitly with respect to an
 377 actual rendering function. Therefore, the decoder of GQN, which is another neural network based
 378 on ConvLSTM, is expected to learn rendering-like operations solely from an objective that aims
 379 to maximize the log-likelihood of each observation given other observations of the same scene as
 380 context. To control for the difference of network architecture, we train DVR using the same image
 381 encoder architecture as GQN, since when we used ResNet18 as an image encoder, GQN did not
 382 converge.

383 Examples of generated images of Geons from GQN are shown in Figure 5 (Left). As we can see,
 384 GQN successfully captures the object from novel viewpoints.

385 To assess the power of representations learned by GQN in the same way as DVR, we take the
 386 representation network and add a linear layer on top. We then finetune the linear layer on 10-Geon
 387 classification, while freezing the rest of the weights. We compare this model to the architecture-
 388 controlled version of the DVR-Last model.

389 Since GQN can take more than one view of images, we prepare 6 models that are finetuned based on
 390 either of {1, 2, 4, 8, 16, 32}-views. The resulting test accuracy of finetuned GQN encoders against
 391 the number of views is shown in Figure 5 (Right). Despite the strong viewpoint generalization of
 392 GQN, we see that finetuned GQN requires more than 2 views (i.e., 3 or 4 views) to reach the DVR
 393 level accuracy, and only outperforms DVR after we feed more than 8 views. This suggests that the
 394 inductive bias from 3D inference is more efficient to obtain good representations.

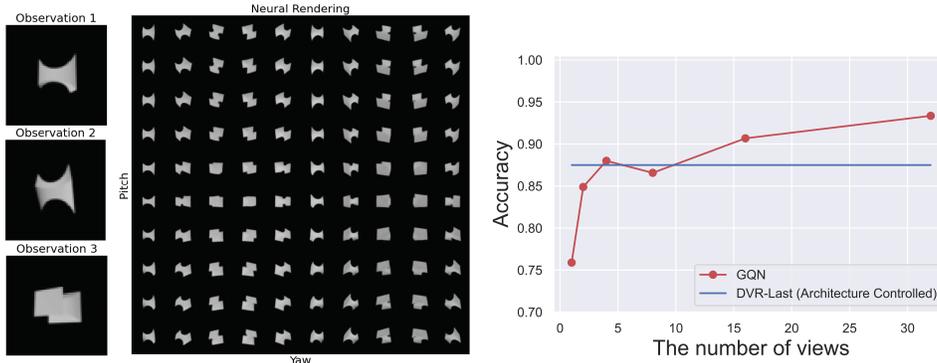


Figure 5: Left: Example Geon images rendered from GQN based on 3 views. Right: GQN Test Accuracy v.s. the number of views. As a reference, we also plot the 1-view DVR accuracy. Here, we used the same architecture for the image encoders of DVR and GQN.

395 B.1 Adversarial Robustness

396 In Figure 6, we provide additional results for adversarial robustness, where we attack AT- L_2 using
 397 L_∞ -PGD. Similar to the case of AT- L_∞ , we see that 3D pretraining improves robustness over the
 398 vanilla AT models for all background settings.

399 B.2 Robustness to Common Corruptions

400 In this section, we provide additional results for common corruptions. In Table 6, we provide the re-
 401 sults for the black background setting. Here again we see that 3D pretraining further improves vanilla
 402 AT models. In Table 7, we provide more detailed results of distributional shift in the backgrounds.
 403 Even after adding image corruptions, we still see that DVR+AT performs best, confirming that 3D
 404 shape bias from 3D pretraining complements the performance of AT to increase model robustness.

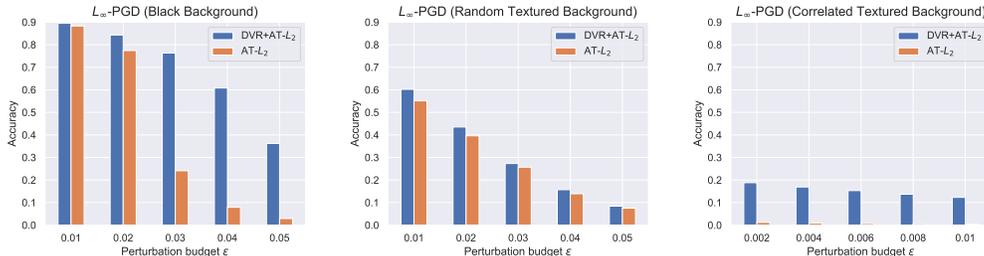


Figure 6: Robustness comparison between AT- L_2 and DVR+AT- L_2 with increasing perturbation budget ϵ on three variations of Geon3D-10. We attack our models using L_∞ -PGD with 100 iterations and $\epsilon/10$ to be the stepsize.

Table 6: Accuracy of shape-biased classifiers against common corruptions under unseen views on Geon3D-10 (black backgrounds).

	REGULAR	INFODROP	STYLIZED	AT- L_2	AT- L_∞	DVR+AT- L_2	DVR+AT- L_∞
INTACT	0.866	0.845	0.822	0.908	0.910	0.912	0.92
PIXELATE	0.685	0.773	0.781	0.905	0.910	0.911	0.919
DEFOCUS BLUR	0.303	0.247	0.755	0.900	0.909	0.897	0.909
GAUSSIAN NOISE	0.548	0.291	0.803	0.620	0.885	0.914	0.919
IMPULSE NOISE	0.140	0.190	0.750	0.542	0.100	0.916	0.918
FROST	0.151	0.323	0.783	0.140	0.100	0.22	0.3
FOG	0.138	0.163	0.764	0.100	0.100	0.119	0.149
ELASTIC	0.612	0.635	0.617	0.628	0.664	0.645	0.655
JPEG	0.799	0.821	0.810	0.905	0.911	0.912	0.92
CONTRAST	0.510	0.180	0.772	0.163	0.258	0.213	0.335
BRIGHTNESS	0.552	0.832	0.818	0.160	0.137	0.385	0.931
ZOOM BLUR	0.475	0.462	0.748	0.891	0.917	0.902	0.92

405 C Related Work and Discussions

406 **3D datasets.** Geon3D is smaller in scale and less complex in shape variation relative to some of the
 407 existing 3D model datasets, including ShapeNet [8] and ModelNet [43]. These datasets have been
 408 instrumental for recent advances in 3D computer vision models (e.g. Niemeyer et al. [28], Sitzmann
 409 et al. [35]). However, at a practical level, these 3D model datasets are not yet suitable for our goal
 410 (which is to establish whether introducing 3D shape bias into vision models induce robustness):
 411 Even though existing learning-based 3D reconstruction models can perform well when trained on
 412 a single or a very small number of categories from these datasets, these models do not scale well
 413 with increasing number of object categories. For example, on ShapeNet, when these models are
 414 required to learn a non-trivial number of object categories (e.g., 10 or more) at the same time, the
 415 resulting 3D shape reconstructions degrade significantly, unable to capture many salient aspects of
 416 shape variation across and within categories. For us, such failure confounds inferences we can make
 417 about the role of shape bias in robustness, which is our central question: Would a negative result be
 418 because the model does not perform well on the reconstruction task to begin with or is it that shape
 419 bias has no benefit for robustness? We deliberately designed Geon3D to allow us to take advantage
 420 of the state-of-the-art in learning-based 3D reconstruction models (in this work, the DVR model): It
 421 provides a non-trivial number of distinct shape categories, with considerable shape variation within
 422 and across categories, yet remain tractable to learn by these existing models. As we demonstrate
 423 in this work, despite its simplicity relative to these larger datasets, Geon3D reveals that the current
 424 vision models struggle with image corruptions and that 3D shape bias induces robustness. Our results
 425 based on Geon3D provide compelling evidence that to achieve robustness against distributional shifts
 426 and adversarial examples, a promising and effective approach is to build models with 3D shape bias.
 427 In future work, we are excited to explore this hypothesis in the context of more complex shapes and
 428 real-world objects and scenes.

429 **Analysis-by-synthesis.** Our proposal of using 3D inference to achieve robust vision shares the
 430 same goal as analysis-by-synthesis [23, 41, 40]. In DVR, we can see its encoder as a recognition

Table 7: Accuracy of shape-biased classifiers against common corruptions under unseen views on Geon3D-10 with textured background swap.

	REGULAR	INFODROP	STYLIZED	AT- L_2	AT- L_∞	DVR+AT- L_2	DVR+AT- L_∞
INTACT	0.045	0.121	0.268	0.015	0.311	0.219	0.439
PIXELATE	0.044	0.096	0.275	0.017	0.306	0.201	0.415
DEFOCUS BLUR	0.044	0.093	0.268	0.024	0.242	0.206	0.338
GAUSSIAN NOISE	0.046	0.160	0.269	0.015	0.320	0.209	0.408
IMPULSE NOISE	0.058	0.096	0.228	0.015	0.078	0.207	0.147
FROST	0.020	0.138	0.255	0.070	0.149	0.144	0.227
FOG	0.032	0.114	0.273	0.077	0.099	0.149	0.124
ELASTIC	0.044	0.109	0.260	0.100	0.196	0.176	0.264
JPEG	0.041	0.089	0.264	0.016	0.306	0.206	0.419
CONTRAST	0.055	0.107	0.274	0.066	0.090	0.148	0.126
BRIGHTNESS	0.036	0.127	0.268	0.026	0.270	0.189	0.379
ZOOM BLUR	0.081	0.082	0.290	0.032	0.269	0.249	0.375

431 network [12], mapping 2D images to their underlying shape, appearance, and pose parameters under
 432 a structured generative model based on a neural rendering function. Even though previous work
 433 considered adversarial robustness of variational autoencoders [33], our study is first to evaluate
 434 robustness arising from analysis-by-synthesis type computations under 3D scenes.

435 D Datasheet

436 A line of work in psychophysics of human visual cognition have argued that the visual system exploits
 437 certain types of shape features in inferring 3D structure and geometry. In Geon3D, by treating these
 438 shape features as the dimensions of variation, we model 40 classes of 3D objects, and render them
 439 from random viewpoints, resulting in an image set and their corresponding camera matrices.

440 **Data Preparation** We construct each Geon using Blender —an open-source 3D computer graphics
 441 software [7].

442 An advantage of Geons over other geometric primitives such as superquadrics [4] is that the shape
 443 categorization of Geons is qualitative rather than quantitative. Thus, each Geon category affords a
 444 high degree of in-class shape deformation, as long as the four defining features of each shape class
 445 remains the same. Such flexibility allows us to construct a number of different 3D model instances
 446 for each Geon class by expanding or shrinking the object along the x, y, or z-axis. For each axis, we
 447 evenly sample the 11 scaling parameters from the interval [0.5, ..., 1.5] with a step size 0.1, resulting
 448 in 1331 3D model instances for each Geon category.

449 **Rendering and data splits** We randomly sample 50 camera positions from a sphere with the object
 450 at the origin. For each model instance, 50 images are rendered using these camera positions with
 451 resolution of 224x224. We then split the data into train/validation/test with ratio 8:1:1 using model
 452 instance ids, where each instance id corresponds to the scaling parameters described above. We also
 453 make sure that all Geon categories are uniformly sampled in each of train/validation/test sets.

454 **Dataset distribution** The full Geon3D-40 (black background) will be available for download after
 455 publication. Geon3D is distributed under the CC BY-SA 4.0 license.¹ We plan to maintain different
 456 versions of Geon3D as we extend the dataset to include more complicated objects by combining
 457 Geon3D as parts. The authors bear all responsibility in case of violation of rights and confirmation
 458 of the data license. Upon publication, the dataset website will become available, where we will add
 459 structured metadata to a dataset’s meta-data page, a persistent dereferenceable identifier, and any
 460 future updates.

461 **How to use Geon3D** Our dataset contains 40 Geon categories, where each folder contains 1331
 462 subfolders. The name of the subfolder represents the scaling factors for the x, y, and z direction. For

¹<https://creativecommons.org/licenses/by-sa/4.0/legalcode>

463 example, 0.5_1.0_1.3 means the Geon model is scaled by 0.5, 1.1, and 1.3 for x, y, and z axis,
 464 respectively. Each subfolder contains the 'rgb' folder, 'mask' folder, and 'pose' folder. The 'rgb'
 465 folder contains 50 images taken from 50 random viewpoints. The 'mask' and 'pose' folders are used
 466 for 3D reconstruction tasks. An example code will be provided to demonstrate how to load these
 467 'mask' and 'pose' information to do 3D reconstruction task.

468 **Benchmarking metric** Our metric for benchmarking model robustness is accuracy under different
 469 noise types (e.g. Section 3.1, 3.2, 3.3, 3.4). Unless we achieve near-perfect accuracy on each noise
 470 type, we don't think robustness issues are solved on this dataset. We would like to avoid using a
 471 single metric such as the mean robust accuracy, since such a metric inevitably obscures the intricate
 472 differences that arise from different noise types.

473 **List of 40 Geons** In Figure 7, we provide a list of 40 Geons we have constructed. The label for each
 474 Geon class represents the four defining shape features, in the order of "axis", "cross section", "sweep
 475 function", "termination", as described in the main paper. We put "na" for the termination when the
 476 sweep function is constant. We also distinguish the two termination types "c-inc" and "c-dec" when
 477 the sweep function is monotonic. For instance, "c-inc" means that the curved surface is at the end
 478 of the increasing sweep function, whereas "c-dec" means that the curved surface is at the end of
 479 the decreasing sweep function. As a reference, here is the mapping between the name and the code
 480 of 10 Geons we used in 10-Geon classification: "Arch": c_s_c_na, "Barrel": s_c_ec_t, "Cone":
 481 s_c_m_p, "Cuboid": s_s_c_na, "Cylinder": s_c_c_na, "Truncated cone": s_c_m_t, "Handle":
 482 c_c_c_na, "Expanded Handle": c_c_m_t, "Horn": c_c_m_p, "Truncated pyramid": s_s_m_t.

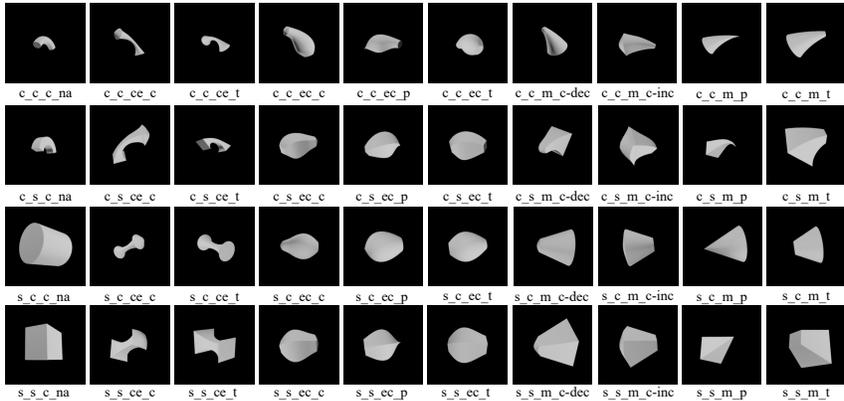


Figure 7: The list of 40 Geons we constructed.

483 E Reproducibility: Training details

484 We used GeForce RTX 2080Ti GPUs for all of our experiments. GQN training takes about a week
 485 until convergence on a single GPU. DVR 3D reconstruction training takes roughly about 1.5 days on
 486 a single GPU. The hyperparameters for 10-Geon classification, described in the main paper, were
 487 chosen by monitoring the model convergence on the validation set. The inputs to all models during
 488 classification are only RGB images. (Camera matrices are only used for the rendering module during
 489 pretraining for 3D reconstruction.)

490 **DVR** We used the code ² open-sourced by Niemeyer et al. [28]. We followed the default hyperpa-
 491 rameters recommended by Niemeyer et al. [28] for 3D reconstruction training, with the exception of
 492 batch size, which we set 32 to fit into a single GPU memory.

493 **Adversarial Training** Through extensive experiments, Zhang & Zhu [42] demonstrate that AT
 494 models develop 2D shape bias, which is considered to explain, in part, the strong adversarial

²https://github.com/autonomousvision/differentiable_volumetric_rendering

495 robustness of AT models. In our experiments, we use L_∞ and L_2 based adversarial training. We used
 496 the python package ³ to perform adversarial training. For $AT(L_2)$, we use attack steps 7, epsilon 3.0,
 497 attack lr 0.5. For $AT(L_\infty)$, we use attack steps 7, epsilon 0.05, attack lr 0.01. use best (final) PGD
 498 step as example. Both models trained for 70 epochs with batch size 100, which was sufficient for
 499 model convergence.

500 **GQN** We used the open-source code ⁴ to implement our GQN. Due to the training instability, we
 501 rescale the image size from 224 x 224 to 64 x 64.

502 **InfoDrop** We used the original author’s implementation ⁵. The method exploits the fact that texture
 503 often repeats itself, and hence is highly correlated with and can be predicted by the texture information
 504 in the neighboring regions, whereas shape-related features such as edges and contours are less coupled
 505 at the locality of neighboring regions.

506 **Stylized** We follow the same protocol as [15] by replacing the texture of each image of Geon3D-10
 507 by a randomly selected texture from paintings through the AdaIn style-transfer algorithm [20]. To
 508 stylize Geon3D, we used the code ⁶ introduced by the original author of Stylized-ImageNet [15].

509 **Dataset** For training Geon3D image classifiers, we center and re-scale the color values of Geon3D
 510 with $\mu = [0.485, 0.456, 0.406]$ and $\sigma = [0.229, 0.224, 0.225]$, which is estimated from ImageNet.
 511 We construct the 40 3D model instances as well as the whole training data in Blender. We then
 512 normalize the object bounding box to a unit cube, which is represented as 1.0_1.0_1.0 in the
 513 dataset folder.

514 **Background textures** We used the following label-to-texture class mapping: {0: 'zigzagged', 1:
 515 'banded', 2: 'wrinkled', 3: 'striped', 4: 'grid', 5: 'polka-dotted', 6: 'chequered', 7: 'blotchy', 8:
 516 'lacelike', 9: 'crystalline' }. For the distributional shift experiment we used the following mapping: {
 517 0: 'crystalline', 1: 'zigzagged', 2: 'banded', 3: 'wrinkled', 4: 'striped', 5: 'grid', 6: 'polka-dotted', 7:
 518 'chequered', 8: 'blotchy', 9: 'lacelike', }. The DTD data is licensed under the Creative Commons
 519 Attribution 4.0 License. ⁷

520 **Evaluation set** For all the evaluation sets in the experiment section, we used the same subset of the
 521 test split, where we randomly pick 1000 model instance ids, and randomly sample 1 view out of 50
 522 views for every model instance.

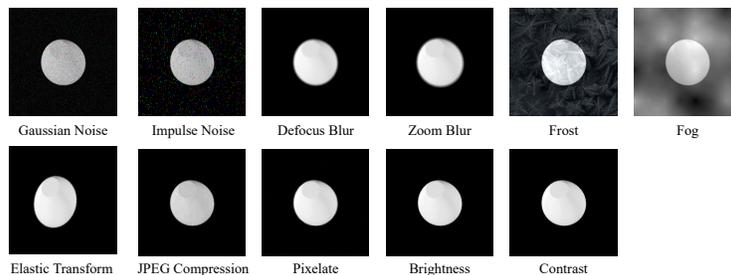


Figure 8: Examples of image corruptions.

523 We use the original author’s code ⁸ to generate common corruptions shown in Figure 8.

³<https://github.com/MadryLab/robustness>

⁴<https://github.com/iShohei220/torch-gqn>

⁵<https://github.com/bfshi/InfoDrop>

⁶<https://github.com/bethgelab/stylize-datasets>

⁷<https://creativecommons.org/licenses/by/4.0/>, <https://www.tensorflow.org/datasets/catalog/dtd>

⁸<https://github.com/hendrycks/robustness>