

---

# Finding Biological Plausibility for Adversarially Robust Features via Metameric Tasks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Recent work suggests that feature constraints in the training datasets of deep neu-  
2       ral networks (DNNs) drive robustness to adversarial noise (Ilyas et al., 2019).  
3       The representations learned by such adversarially robust networks have also been  
4       shown to be more human perceptually-aligned than non-robust networks via image  
5       manipulations (Santurkar et al., 2019; Engstrom et al., 2019). Despite appearing  
6       closer to human visual perception, it is unclear if the constraints in robust DNN  
7       representations match biological constraints found in human vision. Human vision  
8       seems to rely on texture-based/summary statistic representations in the periphery,  
9       which have been shown to explain phenomena such as crowding (Balas et al., 2009)  
10      and performance on visual search tasks (Rosenholtz et al., 2012). To understand  
11      how adversarially robust optimizations/representations compare to human vision,  
12      we performed a psychophysics experiment using a metamer task similar to Freeman  
13      & Simoncelli (2011); Wallis et al. (2019); Deza et al. (2017) where we evaluated  
14      how well human observers could distinguish between images synthesized to match  
15      adversarially robust representations compared to non-robust representations and a  
16      texture synthesis model of peripheral vision (Texforms (Long et al., 2018)). We  
17      found that the discriminability of robust representation and texture model images  
18      decreased to near chance performance as stimuli were presented farther in the  
19      periphery. Moreover, performance on robust and texture-model images showed  
20      similar trends within participants, while performance on non-robust representa-  
21      tions changed minimally across the visual field. These results together suggest  
22      that (1) adversarially robust representations capture peripheral computation better  
23      than non-robust representations and (2) robust representations capture peripheral  
24      computation similar to current state-of-the-art texture peripheral vision models.  
25      More broadly, our findings support the idea that localized texture summary statist-  
26      ic representations may drive human invariance to adversarial perturbations and  
27      that the incorporation of such representations in DNNs could give rise to useful  
28      properties like adversarial robustness.

## 29 1 Introduction

30      Texture-based summary statistic models of the human periphery have been shown to explain key  
31      phenomena such as crowding (Balas et al., 2009; Freeman & Simoncelli, 2011) and performance  
32      on visual search tasks (Rosenholtz et al., 2012) when used to synthesize feature-matching images.  
33      These analysis-by-synthesis models have also been used to explain mid-level visual computation (*e.g.*  
34      V2) via perceptual discrimination tasks on images for humans and primates (Freeman & Simoncelli,  
35      2011; Ziemba et al., 2016; Long et al., 2018).

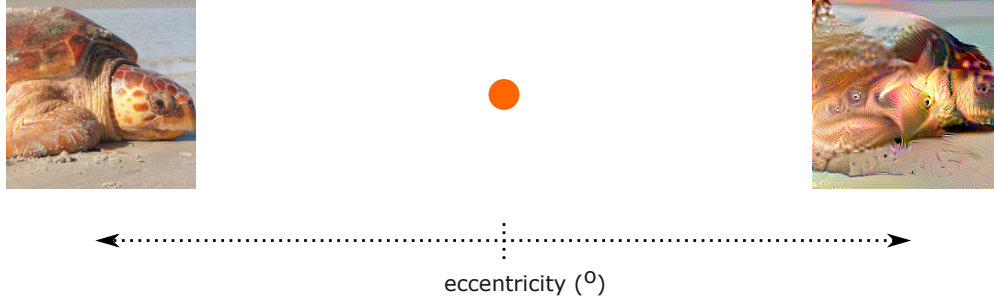


Figure 1: A sample un-perturbed (left) and synthesized adversarially robust (right) image are shown peripherally. When a human observer fixates at the orange dot (center), both images – now placed away from the fovea – are perceptually indistinguishable to each other (i.e. *metameric*). In this paper we investigate if there is a relationship between peripheral representations in humans and learned representations of adversarially trained networks in machines in an analysis-by-synthesis approach. We psychophysically test this phenomena over a variety of images synthesized from an adversarially trained network, a non-adversarially trained network, and a model of peripheral computation as we manipulate retinal eccentricity over 12 humans subjects.

36 While summary statistic models can succeed at explaining peripheral computation in humans, they  
 37 fail to explain foveal computation and core object recognition that involve other representational  
 38 strategies (Logothetis et al., 1995; Riesenhuber & Poggio, 1999; DiCarlo & Cox, 2007; Hinton, 2021).  
 39 Modelling foveal vision with deep learning indeed has been the focus of nearly all object recognition  
 40 systems in computer vision (as machines do not have a periphery) (LeCun et al., 2015; Schmidhuber,  
 41 2015). Despite their success, however, they are vulnerable to adversarial perturbations. This phenom-  
 42 ena indicates: 1) a critical failure of current artificial systems (Goodfellow et al., 2014; Szegedy et al.,  
 43 2013); and 2) a perceptual mis-alignment of such systems with humans (Golan et al., 2019; Feather  
 44 et al., 2019; Firestone, 2020; Geirhos et al., 2021; Funke et al., 2021) – with some exceptions (Elsayed  
 45 et al., 2018). Indeed, there are many strategies to alleviate these sensitivities to perturbations, such as  
 46 data-augmentation (Rebuffi et al., 2021), biologically-plausible inductive biases (Dapello et al., 2020;  
 47 Reddy et al., 2020; Jonnalagadda et al., 2021), and adversarial training (Tsipras et al., 2018; Madry  
 48 et al., 2017). This last strategy in particular (adversarial training) is popular, but has been criticized as  
 49 being non-biologically plausible – despite yielding some perceptually aligned images when inverting  
 50 their representations (Engstrom et al., 2019; Santurkar et al., 2019).

51 We know machines do not have peripheral computation, yet are susceptible to a type of adversarial  
 52 attacks that humans are not. We hypothesize that object representation arising in human peripheral  
 53 computation holds a critical role for high level robust vision in perceptual systems, but testing this has  
 54 not been done. Inspired by recent works that test have tested summary statistic models via metameric  
 55 discrimination tasks (Deza et al., 2017; Wallis et al., 2016, 2017, 2019), we can evaluate how well the  
 56 adversarially robust CNN model approximates the types of computations present in human peripheral  
 57 vision with a set of rigorous psychophysical experiments with respect to synthesized stimuli (Figure 1).  
 58 We evaluated the rates of human perceptual discriminability as a function of retinal eccentricity across  
 59 the synthesized stimuli from an adversarially trained network vs synthesized stimuli from models  
 60 of mid-level/peripheral computation. If the decay rates of perceptual discriminability are similar  
 61 across stimuli, then it suggests that the transformations learned in an adversarially trained network  
 62 are isomorphic to the transformations done by models of peripheral computation – and thus, to the  
 63 human visual system.

## 64 2 Human Psychophysics: Discriminating between stimuli as a function of 65 retinal eccentricity

66 We designed two human psychophysical experiments: the first was a an oddity task similar to  
 67 Wallis et al. (2016), and the second was a matching, two-alternative forced choice task (2AFC).  
 68 Two different tasks were used to evaluate how subjects viewed synthesized images both only in the  
 69 periphery (oddy) and those they saw in the fovea (matching 2AFC). The oddity task consisted of  
 70 finding the oddball stimuli out of a series of 3 stimuli shown peripherally one after the other (100ms)

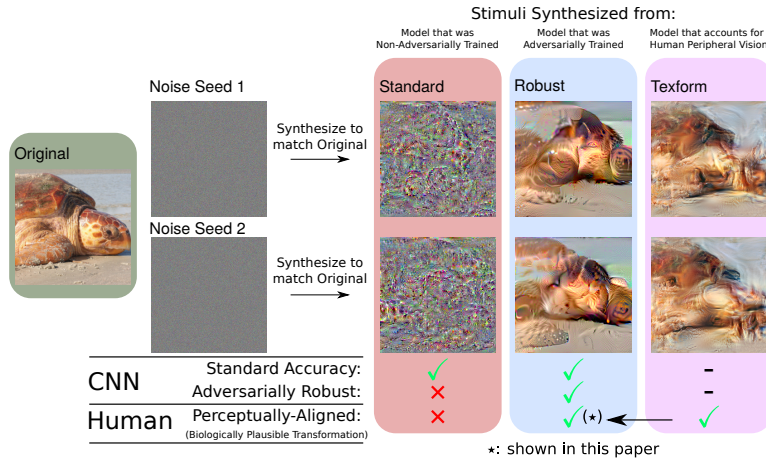


Figure 2: A sub-collection of synthesized stimuli used in our experiments that show differences across (columns) and within (rows) perceptual models. The original stimuli is shown on the left, with two parallel Noise Seeds that give rise to synthesized samples for the Standard, Robust and Texform stimuli. Critically, an adversarially trained network – which was used to synthesize the Robust stimuli (Engstrom et al., 2019) – has implicitly learned to encode a structural prior with localized texture-like distortions similar to the physiologically motivated Texforms that account for several phenomena of *human peripheral computation* (Freeman & Simoncelli, 2011; Rosenholtz et al., 2012; Long et al., 2018). However, Standard stimuli, which are images synthesized from a network with Regular (Non-Adversarial) training have no resemblance to the original sample.

71 masked by empty intervals (500ms) while holding center fixation. Chance for the oddity task was 1  
72 out of 3 (33.3%). The matching 2AFC task consisted of viewing a stimulus in the fovea (100ms) and  
73 then matching it to two candidate templates in the visual periphery (100 ms) while holding fixation.  
74 A 1000 ms mask was used in this experiment and chance was 50%.

75 We used 3 types of stimuli in our experiments: Standard stimuli which were synthesize by a non-  
76 adversarially (standard) trained networks, Robust stimuli which are synthesized by an adversarially  
77 trained stimuli, and Texform stimuli which are synthesized by a model of peripheral and mid-ventral  
78 computation (Figure 2). More information on stimuli synthesis can be seen in Appendix A.

79 For both experiments, we also had interleaved trials where observers had to engage in a Original  
80 stimuli vs Synthesized stimuli task, or a Synthesized stimuli vs Synthesized stimuli discrimination  
81 task (two stimulus pairs synthesized from *different* noise seeds to match model representations). The  
82 goal of these experimental variations (called ‘stimulus roving’) was two-fold: 1) to add difficulty to  
83 the tasks thus reducing the likelihood of ceiling effects; 2) to gather two psychometric functions per  
84 family of stimuli, which portrays a better description of each stimuli’s perceptual signatures.

85 We had 12 participants complete both the oddity and matching 2AFC experiments as shown in  
86 Figure 3. The oddity task was always performed first so that subjects would never have foveated  
87 on the images before seeing them in the periphery. We had two stimulus conditions (1) robust &  
88 standard model images and (2) texforms. Condition 1 consisted of the inverted representations of  
89 the adversarially robust and standard-trained models. The two model representations were randomly  
90 interleaved since they were synthesized with the same procedure. Condition 2 consisted of texforms  
91 synthesized with a fixed and perceptually optimized fixation and scaling factor which yielded images  
92 closest in structure to the robust representations at foveal viewing (robust features have no known  
93 fixation and scaling – which is why we evaluate multiple points in the periphery. Recall Figure 5).  
94 We randomly assigned the order in which participants saw the different stimuli.

95 The main results of our 2 experiments can be found in Figure 4, where we show how well Humans  
96 can discriminate per type of stimuli class and task. Mainly Human observers achieve near perfect  
97 discrimination rates for the Standard stimuli wrt to their original references, but near chance levels  
98 when discriminating to another synthesized sample. This occurs for both experimental paradigms  
99 (Oddity + 2AFC), suggesting that the network responsible for encoding standard stimuli is a poor  
100 model of human peripheral vision given no interaction with retinal eccentricity.

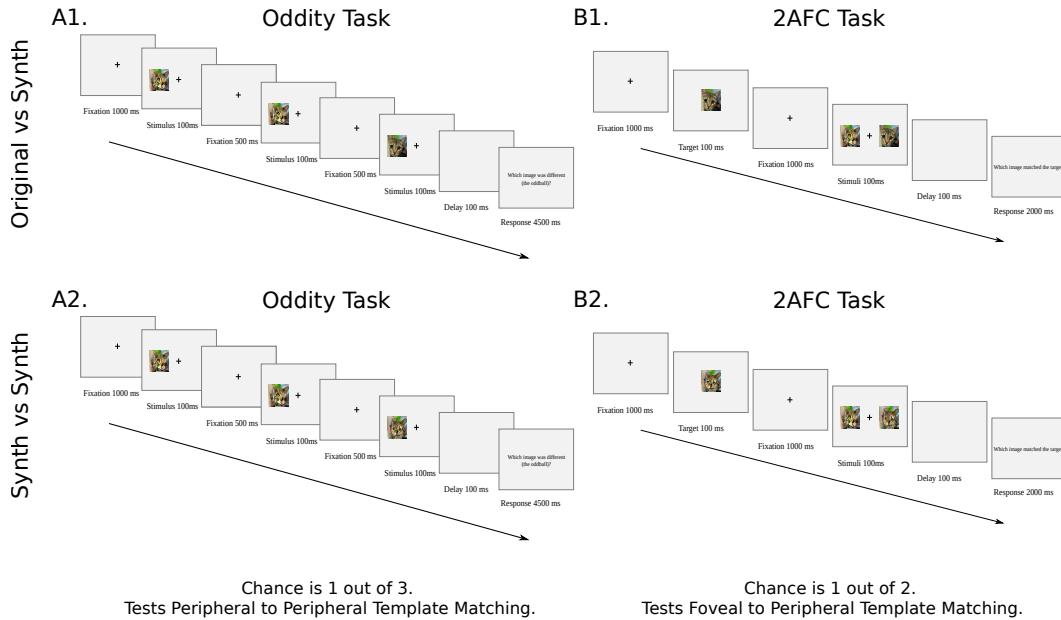


Figure 3: A schematic of the two human psychophysics experiments conducted in our paper. The first (A1.,A2.) illustrates an Oddity task where observers must determine the ‘oddball’ stimuli without moving their eyes for brief presentation times (100 ms) which do not allow for eye-movements or feedback processing. The second experiment (B1.,B2.) shows the 2 Alternative Forced Choice (2AFC) Matching Tasks where observers must match the foveal template to 2 potential candidates on the left or right of the image. All trials are done while observers are instructed to remain fixating at the center of the image. Differences across rows indicate the type of interleaved trials shown to the observers: (1) Original vs Synthesized, and (2) Synthesized vs Synthesized.

101 However, we observe that Humans show similar perceptual discriminability rates for Robust and  
 102 Texform stimuli – and that these vary in a similar way as a function of retinal eccentricity. Indeed, for  
 103 both of these stimuli their perceptual discrimination rates appear to follow a sigmoidal decay-like  
 104 curve when comparing the stimuli to the original, and also between synthesized samples. The  
 105 similarity between the blue and magenta curves from Figure 4 suggests that if the texform stimuli do  
 106 capture some aspect of peripheral computation, then – by transitivity – so do the adversarial stimuli  
 107 which were rendered from an adversarially trained network. These results empirically verify our  
 108 initial hypothesis that adversarially trained networks encode a similar set of transformations as the  
 109 human visual periphery. A superposition of these results in reference to the Robust stimuli for a better  
 110 interpretation can also be seen in Figure 4 (B.).

### 111 3 Discussion

112 We found that stimuli synthesized from an adversarially trained (and thus robust) network are  
 113 metameric to the original stimuli in the further periphery (slightly above 30 deg) for both Oddity and  
 114 2AFC Matching tasks. However, more important than deriving a critical eccentricity for metameric  
 115 guarantees across stimuli in Humans – we found a surprisingly similar pattern of results in terms of  
 116 how perceptual discrimination interacts with retinal eccentricity when comparing the adversarially  
 117 trained network’s robust stimuli with classical models of peripheral computation and V2 encoding  
 118 (mid-level vision) that were used to render the texform stimuli (Freeman & Simoncelli, 2011; Long  
 119 et al., 2018; Ziemba et al., 2016; Ziemba & Simoncelli, 2021). Further, this type of eccentricity-driven  
 120 interaction does not occur for stimuli derived from non-adversarially trained (standard) networks.

121 More generally, now that we found that adversarially trained networks encode a similar class of  
 122 transformations that occur in the visual periphery – how do we reconcile the fact that adversarial  
 123 training is biologically *implausible* in humans? Recall from the work of Ilyas et al. (2019) that per-  
 124 forming *standard training* on robust images yielded similar generalization and adversarial robustness

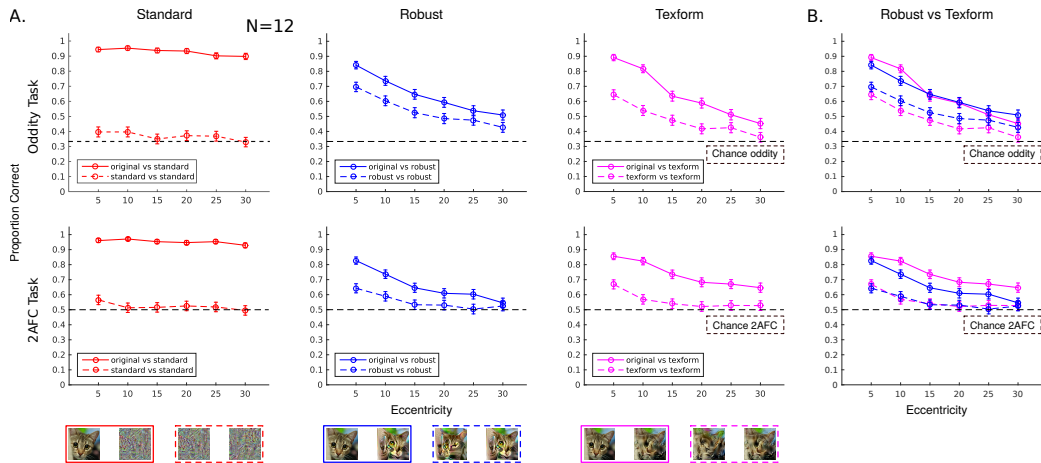


Figure 4: Pooled observer results of both psychophysical experiments are shown (top and bottom row). (A.) Left: we see that observers almost perfectly discriminate the original image from the standard stimuli, in addition to chance performance when comparing against synthesized stimuli. Critically there is no interaction of the standard stimuli with retinal eccentricity which suggests that the model used to synthesize such stimuli is a poor model of peripheral computation. Middle: Human observers do worse at discriminating the robust stimuli from the original as a function of eccentricity and also between synthesized robust samples. Given this decay in perceptual discriminability, it suggests that the adversarially trained model used to synthesize robust stimuli does capture aspects of peripheral computation. This effect is also seen with the texforms (Right) – which have been extensively used as stimuli from derived models that capture peripheral and V2-like computation. (B.) Superimposed human performance for Robust and Texform stimuli. Errorbars are computed via bootstrapping and represent the 95% confidence interval.

125 as performing adversarial training on standard images; how does this connect then to human learning  
 126 if we assume a uniform learning rule in the fovea and the periphery?

127 We think the answer lies in the fact that as humans learn to perform object recognition, they not  
 128 only fixate at the target image, but they also look around, and can eventually learn where to make  
 129 a saccade given candidate object peripheral template – thus learning certain invariances when the  
 130 object is placed both in the fovea and the periphery (Cox et al., 2005; Williams et al., 2008; Poggio  
 131 et al., 2014; Han et al., 2020). This is an idea that dates back to Von Helmholtz (1867) as highlighted  
 132 in Stewart et al. (2020) on the interacting mechanisms of foveal and peripheral vision in humans.

133 Altogether, this could suggest that spatially-uniform high-resolution processing is redundant and  
 134 sub-optimal in the *o.o.d.* regime – as translation invariant adversarially-vulnerable CNNs have  
 135 no foveated/spatially-adaptive computation. Counter-intuitively, the fact that our visual system *is*  
 136 spatially-adaptive could give rise to a more robust encoding mechanism of the visual stimulus as  
 137 observers can encode a distribution rather than a point as they move their center of gaze. Naturally,  
 138 from all the possible types of transformations, the ones that are similar to those shown in this paper –  
 139 which loosely resemble localized texture-computation – are the ones that potentially lead to a robust  
 140 hyper-plane during learning for the observer (See Fig. 9; Appendix).

141 Future work is looking into reproducing the experiments carried out in this paper with a physiological  
 142 component to explore temporal dynamics (MEG) and localization (fMRI) evoked from the stimuli.  
 143 While it is not obvious if we will find a perceptual signature of the adversarial robust stimuli in  
 144 humans, we think this novel stimuli and experimental paradigm presents a first step towards the road  
 145 of linking what is known (and unknown) across texture representation, peripheral computation, and  
 146 adversarial robustness in humans and machines.

## 147 References

- 148 Stuart M Anstis. A chart demonstrating variations in acuity with retinal position. *Vision research*, 14  
149 (7):589–592, 1974.
- 150 Benjamin Balas, Lisa Nakano, and Ruth Rosenholtz. A summary-statistic representation in peripheral  
151 vision explains visual crowding. *Journal of vision*, 9(12):13–13, 2009.
- 152 Alexander Berardino, Valero Laparra, Johannes Ballé, and Eero Simoncelli. Eigen-distortions of  
153 hierarchical representations. *Advances in Neural Information Processing Systems*, 30, 2017.
- 154 Peter J Burt and Edward H Adelson. The laplacian pyramid as a compact image code. In *Readings in*  
155 *computer vision*, pp. 671–679. Elsevier, 1987.
- 156 David D Cox, Philip Meier, Nadja Oertelt, and James J DiCarlo. 'breaking' position-invariant object  
157 recognition. *Nature neuroscience*, 8(9):1145–1147, 2005.
- 158 Joel Dapello, Tiago Marques, Martin Schrimpf, Franziska Geiger, David D Cox, and James J DiCarlo.  
159 Simulating a primary visual cortex at the front of cnns improves robustness to image perturbations.  
160 *BioRxiv*, 2020.
- 161 Arturo Deza and Talia Konkle. Emergent properties of foveated perceptual systems. *arXiv preprint*  
162 *arXiv:2006.07991*, 2020.
- 163 Arturo Deza, Aditya Jonnalagadda, and Miguel Eckstein. Towards metamerism via foveated style  
164 transfer. *arXiv preprint arXiv:1705.10041*, 2017.
- 165 Arturo Deza, Yi-Chia Chen, Bria Long, and Talia Konkle. Accelerated texforms: Alternative methods  
166 for generating unrecognizable object images with preserved mid-level features. In *Conference on*  
167 *Cognitive Computational Neuroscience*, 2019.
- 168 James J DiCarlo and David D Cox. Untangling invariant object recognition. *Trends in cognitive*  
169 *sciences*, 11(8):333–341, 2007.
- 170 Keyan Ding, Kede Ma, Shiqi Wang, and Eero P Simoncelli. Image quality assessment: Unifying  
171 structure and texture similarity. *arXiv preprint arXiv:2004.07728*, 2020.
- 172 Keyan Ding, Kede Ma, Shiqi Wang, and Eero P Simoncelli. Comparison of full-reference image  
173 quality models for optimization of image processing systems. *International Journal of Computer*  
174 *Vision*, 129(4):1258–1281, 2021.
- 175 Gamaleldin F Elsayed, Shreya Shankar, Brian Cheung, Nicolas Papernot, Alexey Kurakin, Ian J  
176 Goodfellow, and Jascha Sohl-Dickstein. Adversarial examples that fool both computer vision and  
177 time-limited humans. In *NeurIPS*, 2018.
- 178 Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Brandon Tran, and Alek-  
179 sander Madry. Adversarial robustness as a prior for learned representations. *arXiv preprint*  
180 *arXiv:1906.00945*, 2019.
- 181 Jenelle Feather, Alex Durango, Ray Gonzalez, and Josh McDermott. Metamers of neural networks  
182 reveal divergence from human perceptual systems. In *Advances in Neural Information Processing*  
183 *Systems*, pp. 10078–10089, 2019.
- 184 Chaz Firestone. Performance vs. competence in human–machine comparisons. *Proceedings of the*  
185 *National Academy of Sciences*, 117(43):26562–26571, 2020.
- 186 Jeremy Freeman and Eero P Simoncelli. Metamers of the ventral stream. *Nature neuroscience*, 14(9):  
187 1195–1201, 2011.
- 188 Christina M Funke, Judy Borowski, Karolina Stosio, Wieland Brendel, Thomas SA Wallis, and  
189 Matthias Bethge. Five points to check when comparing visual perception in humans and machines.  
190 *Journal of Vision*, 21(3):16–16, 2021.
- 191 Leon Gatys, Alexander S Ecker, and Matthias Bethge. Texture synthesis using convolutional neural  
192 networks. *Advances in neural information processing systems*, 28:262–270, 2015.

- 193 Robert Geirhos, Kantharaju Narayanappa, Benjamin Mitzkus, Tizian Thieringer, Matthias Bethge,  
194 Felix A Wichmann, and Wieland Brendel. Partial success in closing the gap between human and  
195 machine vision. *arXiv preprint arXiv:2106.07411*, 2021.
- 196 Wilson S Geisler and Jeffrey S Perry. Real-time foveated multiresolution system for low-bandwidth  
197 video communication. In *Human vision and electronic imaging III*, volume 3299, pp. 294–305.  
198 International Society for Optics and Photonics, 1998.
- 199 Tal Golan, Prashant C Raju, and Nikolaus Kriegeskorte. Controversial stimuli: pitting neural networks  
200 against each other as models of human recognition. *arXiv preprint arXiv:1911.09288*, 2019.
- 201 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial  
202 examples. *arXiv preprint arXiv:1412.6572*, 2014.
- 203 Yena Han, Gemma Roig, Gad Geiger, and Tomaso Poggio. Scale and translation-invariance for novel  
204 objects in human vision. *Scientific reports*, 10(1):1–13, 2020.
- 205 Geoffrey Hinton. How to represent part-whole hierarchies in a neural network. *arXiv preprint*  
206 *arXiv:2102.12627*, 2021.
- 207 Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander  
208 Madry. Adversarial examples are not bugs, they are features. In *Advances in Neural Information*  
209 *Processing Systems*, pp. 125–136, 2019.
- 210 Aditya Jonnalagadda, William Wang, and Miguel P Eckstein. Foveater: Foveated transformer for  
211 image classification. *arXiv preprint arXiv:2105.14173*, 2021.
- 212 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,  
213 2015.
- 214 Gang Liu, Yann Gousseau, and Gui-Song Xia. Texture synthesis through convolutional neural  
215 networks and spectrum constraints. In *2016 23rd International Conference on Pattern Recognition*  
216 *(ICPR)*, pp. 3234–3239. IEEE, 2016.
- 217 Nikos K Logothetis, Jon Pauls, and Tomaso Poggio. Shape representation in the inferior temporal  
218 cortex of monkeys. *Current biology*, 5(5):552–563, 1995.
- 219 Bria Long, Chen-Ping Yu, and Talia Konkle. Mid-level visual features underlie the high-level  
220 categorical organization of the ventral stream. *Proceedings of the National Academy of Sciences*,  
221 115(38):E9015–E9024, 2018.
- 222 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
223 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,  
224 2017.
- 225 Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting  
226 them. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
227 5188–5196, 2015.
- 228 Tomaso Poggio, Jim Mutch, and Leyla Isik. Computational role of eccentricity dependent cortical  
229 magnification. *arXiv preprint arXiv:1406.1770*, 2014.
- 230 Javier Portilla and Eero P Simoncelli. A parametric texture model based on joint statistics of complex  
231 wavelet coefficients. *International journal of computer vision*, 40(1):49–70, 2000.
- 232 Sylvestre-Alvise Rebuffi, Sven Gowal, Dan A Calian, Florian Stimberg, Olivia Wiles, and Tim-  
233 othy Mann. Fixing data augmentation to improve adversarial robustness. *arXiv preprint*  
234 *arXiv:2103.01946*, 2021.
- 235 Manish V Reddy, Andrzej Banburski, Nishka Pant, and Tomaso Poggio. Biologically inspired  
236 mechanisms for adversarial robustness. *arXiv preprint arXiv:2006.16427*, 2020.
- 237 Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object recognition in cortex.  
238 *Nature neuroscience*, 2(11):1019–1025, 1999.

- 239 Ruth Rosenholtz, Jie Huang, Alvin Raj, Benjamin J Balas, and Livia Ilie. A summary statistic  
240 representation in peripheral vision explains visual search. *Journal of vision*, 12(4):14–14, 2012.
- 241 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,  
242 Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition  
243 challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- 244 Shibani Santurkar, Dimitris Tsipras, Brandon Tran, Andrew Ilyas, Logan Engstrom, and Aleksander  
245 Madry. Image synthesis with a single (robust) classifier. *arXiv preprint arXiv:1906.09453*, 2019.
- 246 Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural networks*, 61:85–117,  
247 2015.
- 248 Emma EM Stewart, Matteo Valsecchi, and Alexander C Schütz. A review of interactions between  
249 peripheral and foveal vision. *Journal of vision*, 20(12):2–2, 2020.
- 250 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow,  
251 and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- 252 Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry.  
253 Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- 254 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine  
255 learning research*, 9(11), 2008.
- 256 Hermann Von Helmholtz. *Handbuch der physiologischen Optik: mit 213 in den Text eingedruckten  
257 Holzschnitten und 11 Tafeln*, volume 9. Voss, 1867.
- 258 Thomas SA Wallis, Matthias Bethge, and Felix A Wichmann. Testing models of peripheral encoding  
259 using metamerism in an oddity paradigm. *Journal of vision*, 16(2):4–4, 2016.
- 260 Thomas SA Wallis, Christina M Funke, Alexander S Ecker, Leon A Gatys, Felix A Wichmann, and  
261 Matthias Bethge. A parametric texture model based on deep convolutional features closely matches  
262 texture appearance for humans. *Journal of vision*, 17(12):5–5, 2017.
- 263 Thomas SA Wallis, Christina M Funke, Alexander S Ecker, Leon A Gatys, Felix A Wichmann, and  
264 Matthias Bethge. Image content is more important than bouma’s law for scene metamers. *ELife*, 8:  
265 e42512, 2019.
- 266 Mark A Williams, Chris I Baker, Hans P Op De Beeck, Won Mok Shim, Sabin Dang, Christina  
267 Triantafyllou, and Nancy Kanwisher. Feedback of visual object information to foveal retinotopic  
268 cortex. *Nature neuroscience*, 11(12):1439–1445, 2008.
- 269 Corey M Ziemba and Eero P Simoncelli. Opposing effects of selectivity and invariance in peripheral  
270 vision. *Nature Communications*, 12(1):1–11, 2021.
- 271 Corey M Ziemba, Jeremy Freeman, J Anthony Movshon, and Eero P Simoncelli. Selectivity and  
272 tolerance for visual texture in macaque v2. *Proceedings of the National Academy of Sciences*, 113  
273 (22):E3140–E3149, 2016.



274 **A Synthesizing Stimuli as a window to Model Representation**

275 Suppose we have the functions  $g_{Adv}(\circ)$  and  $g_{Standard}(\circ)$  that represent the adversarially trained and  
 276 standard (non-adversarially) trained neural networks; how can we compare them to human peripheral  
 277 computation if the function  $g_{Human}(\circ)$  is computationally intractable?

278 One solution is to take an analysis-by-synthesis approach and to synthesize a collection of stimuli  
 279 that match the feature response of the model we'd like to analyze – this is also known as feature  
 280 inversion (Mahendran & Vedaldi, 2015; Feather et al., 2019). If the inverted features (stimuli) of two  
 281 models are perceptually similar, then it is likely that the learned representations are also aligned. For  
 282 example, if we'd like to know what is the stimuli  $x'$  that produces the same response to the stimuli  $x$   
 283 for a network  $g'(\circ)$ , we can perform the following minimization:

$$x' =_{x_0} [||g'(x) - g'(x_0)||_2] \tag{1}$$

284 In doing so, we find  $x'$  which should be different from  $x$  for a non-trivial solution. This is known as a  
 285 metameric constraint for the stimuli pair  $\{x, x_0\}$  wrt to the model  $g'(\circ) : g'(x) = g'(x')$  s.t.  $x \neq x'$   
 286 for a starting pre-image  $x_0$  that is usually white noise in the iterative minimization of Eq.1. Indeed,  
 287 for the adversarially trained network of Ilyas et al. (2019); Engstrom et al. (2019); Santurkar et al.  
 288 (2019), we can synthesize robust stimuli wrt to the original image  $x$  via:

$$\tilde{x} =_{x_0} [||g_{Adv}(x) - g_{Adv}(x_0)||_2] \tag{2}$$

289 which implies – if the minimization goes to zero – that:

$$||g_{Adv}(x) - g_{Adv}(\tilde{x})||_2 = 0 \tag{3}$$

290 Recalling the goal of this paper, we'd like to investigate if the following statement is true: “a transfor-  
 291 mation resembling peripheral computation in the human visual system can closely be approximated  
 292 by an adversarially trained network”, which is formally translated as:  $g_{Adv} \sim g_{Human}^{r_*}$  for some  
 293 retinal eccentricity ( $r_*$ ), then from Eq. 3 we can also derive:

$$||g_{Human}^{r_*}(x) - g_{Human}^{r_*}(\tilde{x})||_2 = 0 \tag{4}$$

294 However,  $g_{Human}(\circ)$  is computationally intractable, so how can we compute Eq.4? A first step is  
 295 to perform a psychophysical experiment such that we find a retinal eccentricity  $r_*$  at which human  
 296 observers can not distinguish between the original and synthesized stimuli – thus behaviourally  
 297 proving that the condition above holds, without the need to directly compute  $g_{Human}$ .

298 More generally, we'd like to compare the *psychometric functions* between stimuli generated from a  
 299 standard trained network (standard stimuli), an adversarially trained network (robust stimuli), and  
 300 a model that captures peripheral and mid-level visual computation (textform stimuli (Freeman &  
 301 Simoncelli, 2011; Long et al., 2018)). Then we will assess how the psychometric functions vary as a  
 302 function of retinal eccentricity. If there is significant overlap between psychometric functions between  
 303 one model wrt the model of peripheral computation; then this would suggest that the transformations  
 304 developed by such model are similar to those of human peripheral computation. We predict that this  
 305 will be the case for the adversarially trained network ( $g_{Adv}(\circ)$ ). Formally, for any model  $g$ , and its  
 306 synthesized stimuli  $x_g$  – as shown in Figure 2, we will define the psychometric function  $\delta_{Human}$ ,  
 307 which depends on the eccentricity  $r$  as:

$$\delta_{Human}(g; r) = ||g_{Human}^r(x) - g_{Human}^r(x_g)||_2 \tag{5}$$

308 where we hope to find:

$$\delta_{Human}(g_{Adv}; r) = \delta_{Human}(g_{Textform}; r); \forall r. \tag{6}$$

309 **A.1 Standard and Robust Model Stimuli**

310 To evaluate robust vs non-robust feature representations, we used the ResNet-50 models of Santurkar  
 311 et al. (2019); Ilyas et al. (2019); Engstrom et al. (2019). We used their models so that our results  
 312 could be interpreted in the context of their findings that features may drive robustness. Both models  
 313 were trained on a subset of ImageNet (Russakovsky et al., 2015), termed Restricted ImageNet (Table  
 314 1). The benefit of Restricted ImageNet, stated by Ilyas et al.; Engstrom et al., is models can achieve  
 315 better standard accuracy than on all of ImageNet. One drawback is that it is imbalanced across classes.  
 316 Although the class imbalance was not problematic for comparing the adversarially robust model to

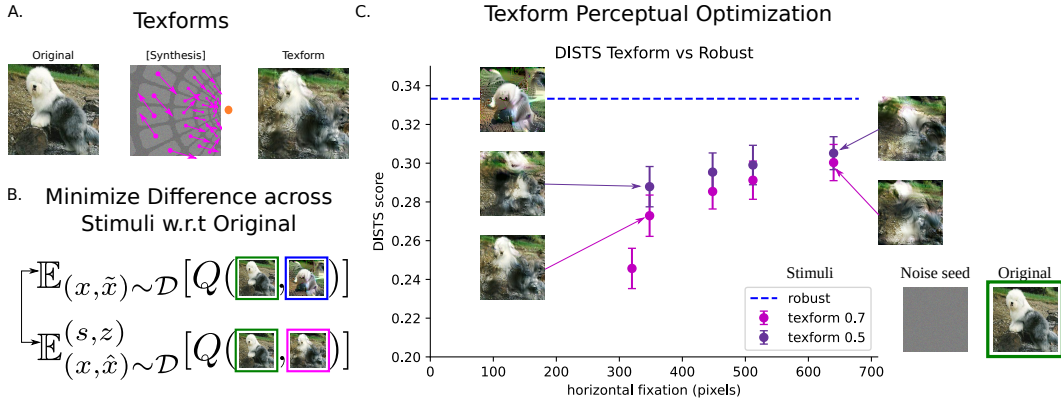


Figure 5: (A.) A cartoon depicting the texform generating process where log-polar receptive fields are used as areas over which localized texture synthesis is performed – imitating the type of texture-based computation found in the human periphery and area V2. (B.) The perceptual optimization framework where the goal is to find the set of texform parameters ( $s_*$ ,  $z_*$ ) over which the loss is minimized to match the levels of distortions of the robust stimuli *before* performing human psychophysics. (C.) The texform perceptual optimization pipeline results show the DISTS scores (Ding et al., 2020) of texforms synthesized across different scaling factors and fixations points compared to adversarially robust stimuli synthesized from the same noise seed across 45 images (5 per RestrictedImageNet class selected randomly). Error bars indicate two standard errors from the mean.

317 standard-trained one, we did ensure that there was a nearly equal number of images per class when  
 318 selecting images for our stimulus set to avoid class effects in our experiment (i.e. people are better at  
 319 discriminating dog examples than fishes independent of the model training).

320 Using their readily available models, we synthesized robust and standard model stimuli using an  
 321 image inversion procedure (Mahendran & Vedaldi, 2015; Gatys et al., 2015; Santurkar et al., 2019;  
 322 Engstrom et al., 2019; Ilyas et al., 2019). We used gradient descent to minimize the difference  
 323 between the representation of the second-to-last network layer of a target image and an initial noise  
 324 seed as shown in Figure 6. Target images were randomly chosen from the test set of Restricted  
 325 ImageNet. We chose 100 target images for each of the 9 classes and synthesized a robust and standard  
 326 stimulus for 2 different noise seeds. 5 target images were later removed as they were gray-scale and  
 327 could not also be rendered as Texforms with the same procedure as the majority. All stimuli were  
 328 synthesized at a size of 256 pixels, this was equivalent to  $7 \times 7$  degree of visual angle (d.v.a.)  
 329 when performing the psychophysical experiments.

## 330 A.2 Texform Stimuli

331 Texforms (Long et al., 2018) are object-equivalent rendered stimuli from the Freeman & Simoncelli  
 332 (2011); Rosenholtz et al. (2012) models that break the metamer constraint to test for mid-level  
 333 visual representations in Humans. These stimuli – initially inspired by the experiments of Balas et al.  
 334 (2009) – preserve the coarse global structure of the image and its localized texture statistics (Portilla  
 335 & Simoncelli, 2000). Critically, we use the texform stimuli – *voiding the metamer constraint* – as a  
 336 perceptual control for the robust stimuli, as the texforms incarnate a sub-class of biologically-plausible  
 337 distortions that loosely resemble the mechanisms of human peripheral processing.

338 As the texform model has 2 main parameters which are the scaling factor  $s$  and the simulated point  
 339 of fixation  $z$ , we must perform a perceptual optimization procedure to find the set of texforms  $\hat{x}$   
 340 that match the robust stimuli  $\tilde{x}$  as close as possible (w.r.t to the original image) *before* testing their  
 341 discriminability to human observers as a function of eccentricity. To do this, we used the accelerated  
 342 texform implementation of Deza et al. (2019) and generated 45 texforms with the *same* collection of  
 343 initial noise seeds as the robust stimuli to be used as perceptual controls. Similar to Deza & Konkle  
 344 (2020) we minimize the perceptual dissimilarity  $\mathcal{Z}$  to find ( $s_*$ ,  $z_*$ ) over this subset of images that we  
 345 will later use in the human psychophysics ( $\sim 900$  texforms):

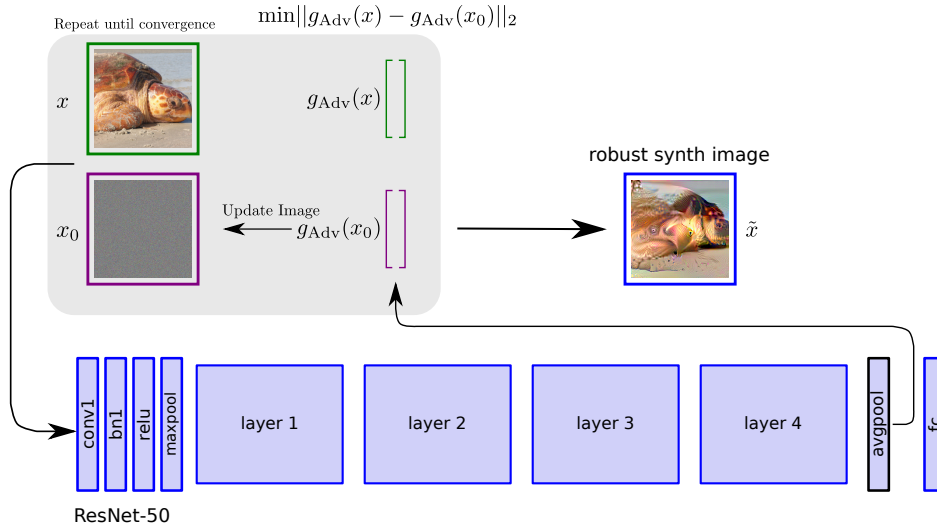


Figure 6: The Robust Image Synthesis pipeline: A noise image  $x_0$  is passed through an adversarially trained ResNet-50 and the penultimate layer features  $g_{Adv}(x_0)$  are matched wrt the original images’ penultimate feature activation  $g_{Adv}(x)$  via an L2 loss, and is repeated until convergence (Santurkar et al., 2019; Engstrom et al., 2019). Critically we use  $g_{Adv}(\circ)$  as a summary statistic of peripheral processing in our experiments.

$$(s_*, z_*) =_{(s,z)} \mathcal{Z} = \|\mathbb{E}_{(x,\hat{x}) \sim \mathcal{D}}[Q(x, \hat{x})] - \mathbb{E}_{(x,\hat{x}) \sim \mathcal{D}}^{(s,z)}[Q(x, \hat{x})]\|_2 \quad (7)$$

346 for an image quality assessment (IQA) function  $Q(\circ, \circ)$ . We selected DISTS in our perceptual  
 347 optimization setup given that it is the IQA metric that is most tolerant to texture-based transforma-  
 348 tions (Ding et al., 2020, 2021). A cartoon illustrating the texform rendering procedure, the perceptual  
 349 optimization framework and the respective results can be seen in Figure 5. In our final experiments  
 350 we used texforms rendered with a simulated scale of 0.5 and horizontal simulated point of fixation  
 351 placed at 640 pixels. Critically, this value is *immutable* and texforms (like robust stimuli) will not vary  
 352 as a function of eccentricity to provide a fair discriminability control in the human psychophysics.  
 353 For a further discussion on texforms and their biological plausibility and/or synthesis procedure,  
 354 please see Supplement B.2.

## 355 B Image Synthesis Details

Classes									
RIN	Dog	Cat	Frog	Turtle	Bird	Primate	Fish	Crab	Insect
IN	151-268	281-285	30-32	33-37	68-100	365-382	389-397	118-121	300-319

Table 1: Classes of RestrictedImageNet (RIN) and the corresponding ImageNet (IN) class ranges.

### 356 B.1 Standard and Robust Stimuli

357 We used the publicly available code from Santurkar et al. (2019); Engstrom et al. (2019); Ilyas  
 358 et al. (2019) found here to synthesize both standard and robust stimuli which were derived from a  
 359 regularly and adversarially trained model respectively: [https://github.com/MadryLab/robust\\_](https://github.com/MadryLab/robust_)  
 360 [representations](https://github.com/MadryLab/robust_)

361 A schematic that illustrates the robust stimuli rendering pipeline can be seen in Figure 6. Standard  
 362 stimuli is generated with the same procedure, and number of iterations, but the network  $g_{Adv}(\circ)$  is  
 363 replaced with  $g_{Standard}(\circ)$  instead.

## 364 B.2 Texform Stimuli

365 Texform stimuli were synthesized using the publicly available code of Deza et al. (2019): <https://github.com/ArturoDeza/Fast-Textforms>

367 The following images (class:[image id's]) were removed as they did not converge:

- 368 • texform0: 0:[49],1:[9],2:[],3:[44],4:[],5:[],6:[10],7:[40],8:[].
- 369 • texform1: 0:[49],1:[9,44],2:[],3:[44],4:[],5:[],6:[10],7:[40],8:[]

370 In addition the following image id's were removed from our psychophysical analysis from the texform  
371 stimuli as they converged to the *exact* same image even when starting from different noise seeds.  
372 This was found while doing a post-hoc IQA analysis as the one shown in Figure ???. These stimuli  
373 only occurred for classes 0 (dog) and 1 (cat):

- 374 • texform: 0:[22,25,26,27,29,93,94,95,96,97,98,99],1:[20,21,22,23,73,74]

375 We found that Standard and Robust stimuli did not have this identical convergence problem over the  
376 900 rendered pairs (1800 stimuli in total for Standard and 1800 in total for Robust).

377 **Note 1a:** A common mis-conception is that Freeman & Simoncelli (2011)-derived stimuli (such  
378 as texforms) *do not* contain structural priors and only performs localized texture synthesis over  
379 smoothly overlapping log-polar receptive fields. This has been investigated with great detail in Wallis  
380 et al. (2016, 2017); Liu et al. (2016) that showed that without spectral constraints it is impossible to  
381 generate metameric images from non-stationary textures for the human observer when showing such  
382 stimuli in the visual periphery. For texforms the metameric constraint is purposely broken because  
383 we'd like to test how a specific biologically-plausible family of transformations (embodied through  
384 the synthesis procedure) interacts with eccentricity when the eccentricity-dependent and scaling  
385 factors texform parameters are fixed. See  $(z_*, s_*)$  from Eq. 7.

386 **Note 1b:** The Freeman & Simoncelli (2011) synthesis model is not equivalent to the Portilla &  
387 Simoncelli (2000) synthesis model. The Freeman & Simoncelli (2011) is a super-ordinate synthesis  
388 model class that locally uses the Portilla & Simoncelli (2000) synthesis model over smoothly  
389 overlapping receptive fields in addition to adding a global structural prior. Texforms are rendered  
390 with the Freeman & Simoncelli (2011) model, by placing the simulated point of fixation *outside* the  
391 image (Long et al., 2018; Deza et al., 2019).

392 **Note 1c:** Usual texform rendering time is about 1 day per image, though the rendering procedure  
393 has been accelerated to the order of minutes as shown in Deza et al. (2019). We used their publicly  
394 available code in our experiments. Thus, it is worth noting that synthesizing texforms in the order of  
395 hundreds of thousands (or millions) for supervised learning experiments – has not been done before  
396 and is computationally expensive (may take months), which is why Figure 2 displays no information  
397 on texform-trained CNN's. This direction is current work.

398 **Note 2:** A first naive criticism to the selection of making texforms fixed and not varying as a function  
399 of eccentricity – given the model they were based on (Freeman & Simoncelli, 2011) – is that they  
400 will not create metameric stimuli. Our anticipated reply to this is three-fold, and partially aligned  
401 with the motivation of Long et al. (2018):

- 402 1. Our goal is *not* to make metameric stimuli out of texforms or robust stimuli, but to examine  
403 how perceptual discriminability rates of a *fixed stimuli* change as a function of retinal  
404 eccentricity. By checking if these perceptual decays are similar (which we show) we can  
405 connect both functions that give rise to these apparently un-related transformations (the  
406 stimuli). Recall Eq. 6.
- 407 2. Having a “metameric texform” that changes as a function of eccentricity would defeat the  
408 purpose of using it as a control in our experiments. Had this been the road taken, we would  
409 now have a control curve that will presumably be horizontal and at chance, providing no  
410 information about how the transformation that gives rise to the robust stimuli is linked to the  
411 texform transformation.
- 412 3. The goal of this paper is *not* to make a foveated metamer model that fools human observers  
413 similar to that of Freeman & Simoncelli (2011); Rosenholtz et al. (2012); Deza et al. (2017);

414  
415  
416

Wallis et al. (2019) that would be based on a foveated adversarially trained network. The previous idea however is highly interesting and is being explored in current work, and this work provides a proof of concept that it is tractable.

### 417 B.3 Sample Stimuli

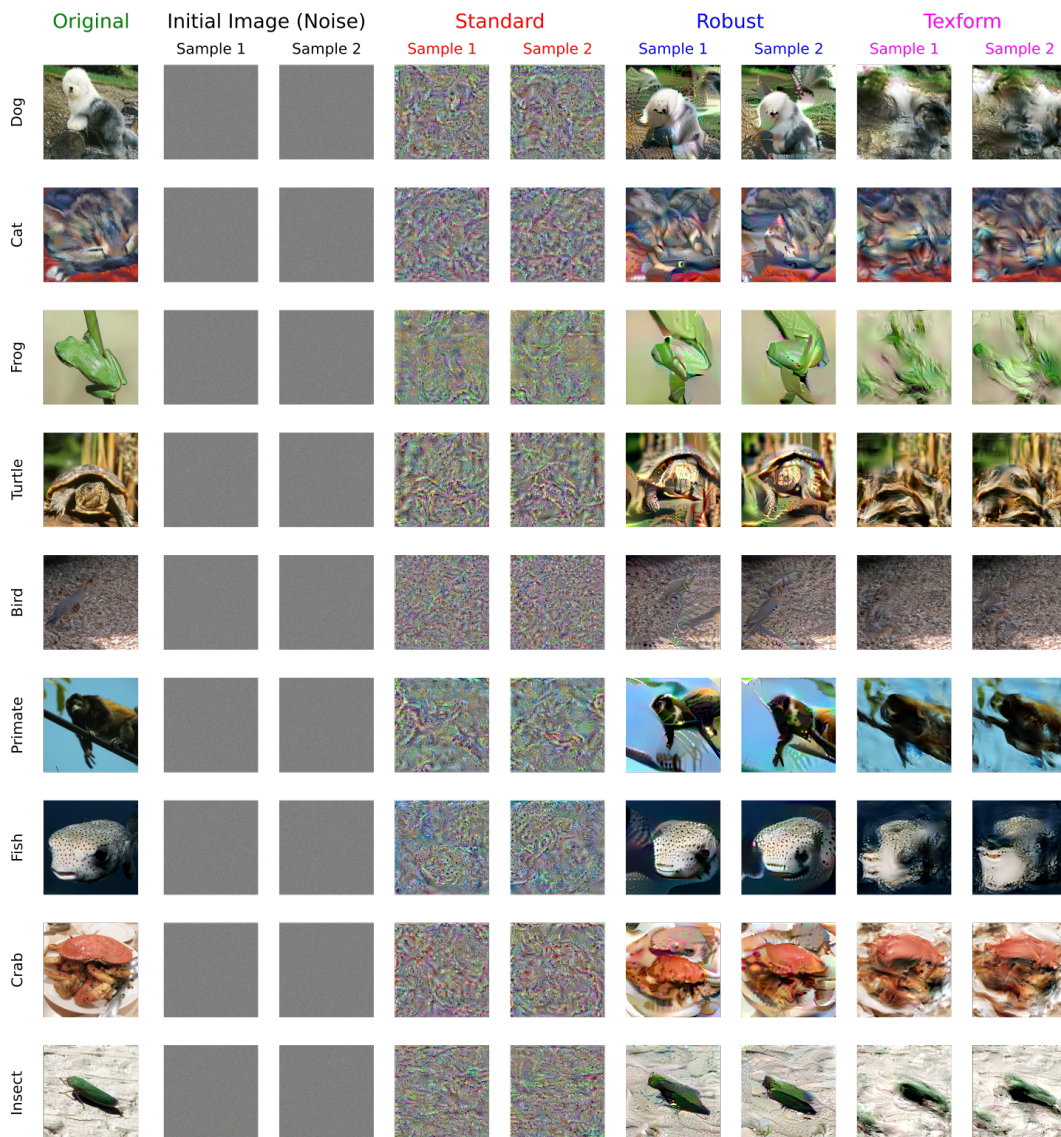


Figure 7: A collection of sample stimuli for each image class used in our experiments.

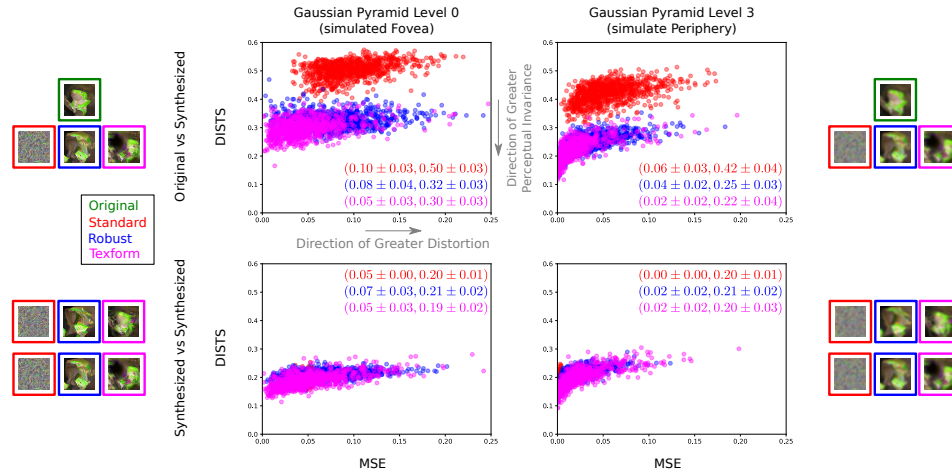


Figure 8: Here we evaluate how the different stimuli differ to each other wrt to the original (top row) or synthesized samples (bottom row) via two IQA metrics: DISTs and MSE. This characterization allows us to compare which model discards more information (MSE) while yielding a greater degree of model based perceptual invariance. We find that Texform and Robust stimuli are similar terms of both IQA scores, suggesting their models compute the same transformations. This is observed at the 0th level (simulated fovea) and 3rd level (simulated periphery) of the Gaussian Pyramid.

### 418 C Simulated Fovea/Periphery Image Quality Assessment (IQA) across 419 stimuli

420 Some distortions are more perceptually noticeable than others for human observers and deep neural  
421 networks (Berardino et al., 2017) – so how do we assess which model better accounts for peripheral  
422 computation, if there are many distortions (derived from the synthesized model stimuli) that can  
423 potentially yield the same perceptual sensitivity in a discrimination task?

424 Our approach consists of computing two IQA metrics (DISTs & MSE) over the entire psychophysical  
425 testing set over 2 opposite levels of a Gaussian Pyramid decomposition (Burt & Adelson, 1987).  
426 This procedure checks which stimuli presents the greatest distortion (MSE), and yet yields greater  
427 perceptual invariance (DISTs). A Gaussian Pyramid decomposition was selected as it stimulates the  
428 frequencies preserved given changes in human contrast sensitivity and cortical magnification factor  
429 from fovea to periphery (Anstis, 1974; Geisler & Perry, 1998). These two metrics were one that is  
430 texture-tolerant and perceptually aligned (DISTs), and another that is a non-perceptually aligned  
431 metric: Mean Square Error (MSE). Both IQA metrics were computed in pixel space for both the  
432 Original vs Synthesized and Synthesized vs Synthesized conditions.

433 Results are explained in Figure 8, where Standard Stimuli yields low perceptual invariance to the  
434 original image at both levels of the Gaussian Pyramid, but robust and texform stimuli have a similar  
435 degree of perceptual invariance. Critically, robust stimuli are slightly more distorted via MSE than  
436 texform stimuli suggesting that the adversarially trained model has learned to represent peripheral  
437 computation better than the texform model by maximizing the perceptual null space and throwing  
438 away more useless low-level image features (hence achieving greater Mean Square Error).

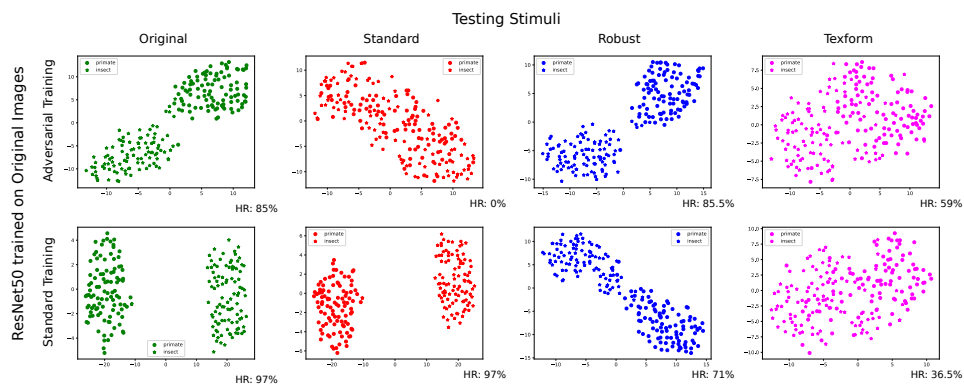


Figure 9: Here we show a 2D projection using t-SNE (Van der Maaten & Hinton, 2008) to visualize the outputs of the last layer of the Adversarially trained network (that was used to synthesize the Robust Stimuli), and the Standard trained network (that was used to synthesize the Standard stimuli), both on a family of different stimuli: Original, Standard, Robust and Texform. The Adversarially trained network – similar to the human – can not distinguish between 2-class Standard Stimuli (unlike the Standard Network that has a near perfect 2-class hit rate). Most importantly, the Adversarially trained network yields a near double hit rate on Texform classification wrt the Standard trained network. This suggests that the Adversarially trained network has a representation that is more perceptually aligned to models of Peripheral Computation than the Standard trained model.