
Small-scale adversarial perturbations expose differences between predictive encoding models of human fMRI responses

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Artificial neural network-based vision encoding models have made significant
2 strides in predicting neural responses and providing insights into visual cognition.
3 However, progress appears to be slowing, with many encoding models achieving
4 similar levels of accuracy in predicting brain activity. In this study, we show that
5 encoding models of human fMRI responses are highly vulnerable to small-scale
6 adversarial attacks, revealing differences not captured using predictive accuracy
7 alone. We then test adversarial sensitivity as a complementary evaluation measure
8 and show that it offers a more effective way to distinguish between highly predictive
9 encoding models. While explicit adversarial training can increase robustness of
10 encoding models, we find that it comes at the cost of brain prediction accuracy.
11 Our preliminary findings also indicate that the choice of model features-to-brain
12 mapping might play a role in optimizing both robustness and accuracy, with sparse
13 mappings typically resulting in more robust encoding models of neural activity.
14 These findings reveal key vulnerabilities of current models, introduce a novel
15 evaluation procedure, and offer a path toward improving the balance between
16 robustness and predictive accuracy for future encoding models.

17 1 Introduction

18 Artificial neural networks (ANNs), loosely inspired by the architecture of the visual cortex, have
19 become the leading models for understanding human vision [1–3]. These models excel not only
20 at complex tasks like object recognition (e.g., ImageNet classification) but also provide a valuable
21 framework for studying visual cognition more broadly [4–6]. ANN-based encoding models, which
22 map neural network features to brain activity, have unlocked a key ability to predict responses at
23 the level of single neurons [7], voxels [8], entire brain regions [9, 10], and even human and non-
24 human primate behavior [11–13]. Early work established a link between a model’s performance on
25 complex tasks (like ImageNet) and the ability to predict brain responses: better task performance
26 typically translated to better brain/behavioral predictions [14, 1, 15]. However, this relationship has
27 plateaued; despite continual improvements in task performance, gains in brain prediction accuracy
28 (henceforth predictivity) have largely stalled. This observation raises critical questions: Are models
29 with similar predictivity learning the same features, or are key differences going unnoticed? Is there a
30 more effective metric that can reveal these differences and help us identify the better models, even
31 when their predictivity appears to be equally high? In this work, we show that small, imperceptible
32 (to humans) adversarial attacks on predictive encoding models can reveal meaningful differences,
33 providing a sharper lens to evaluate their fidelity as models of the brain.

34 The concern that encoding model predictivity has plateaued is not new [14, 9, 10, 15, 16]. This
35 stagnation has sparked two major responses within the field. On one front, this challenge has driven

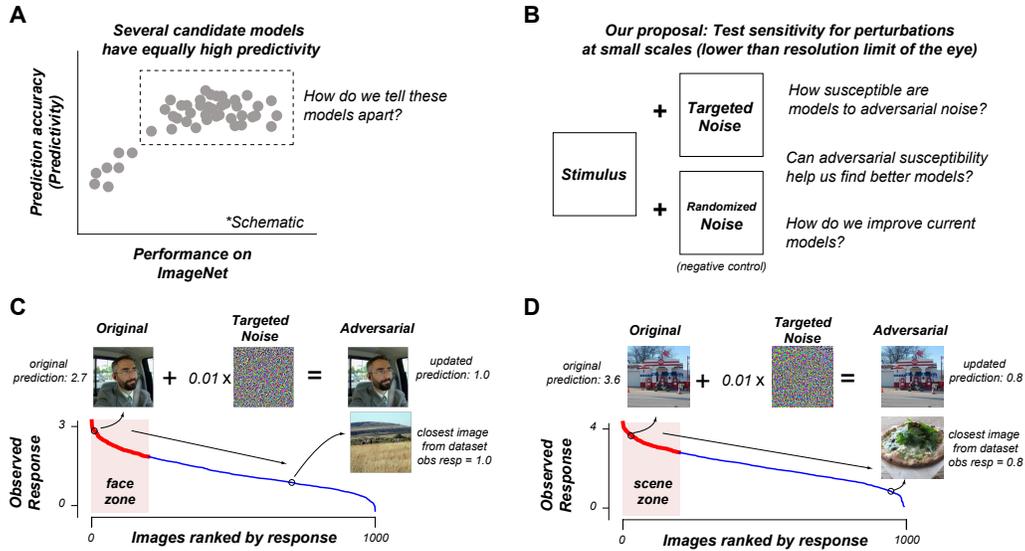


Figure 1: **Motivation, central questions and example adversarial perturbations on encoding models.** **A.** Schematic illustrating the trend observed in previous studies: many encoding models show similarly high predictions on brain data. Performance on ImageNet is shown on the x-axis, and prediction accuracy on brain data is shown on the y-axis. A similar figure with actual data can be found in previous work [1]. **B.** Strategy and central questions. For a given stimulus, we generate a targeted noise pattern (and a randomized control) to assess how sensitive the encoding model is to adversarial noise. **C.** Example of an adversarial perturbation applied to the FFA encoding model (VGG16). The model’s prediction for the original face stimulus (top left) is significantly altered when a targeted imperceptible noise pattern is added (top middle). The modified image produces a much lower response (top right), similar to that for a non-preferred stimulus (bottom). The x-axis shows images sorted by response, and the y-axis represents the FFA’s response. The red region highlights the response for the preferred category (faces). **D.** Same as C, but for an example scene stimulus and the PPA.

36 the development of *entirely new models* incorporating aspects of brain-like operations (like recurrence
 37 [17, 18, 2, 19, 20]) or by directly aligning with behavioral or neural data [21–23]. The second
 38 front challenges predictivity as the primary metric altogether, advocating for *alternative evaluation*
 39 *methods* like centered kernel alignment [24–27] or single-neuron selectivity [28, 29] to capture more
 40 nuanced aspects of brain-model alignment. In this work, we are advocating a slightly different
 41 strategy. Predictivity must remain a vital benchmark measure of our models: predictive models
 42 have enabled new understanding of brain function including the ability to modulate responses in the
 43 visual cortex [30–33]. However, predictivity alone is insufficient, especially when we are limited by
 44 data. We propose complementing predictivity measures with additional evaluation metrics. Here we
 45 introduce adversarial sensitivity as a potential tool for stronger model evaluations.

46 Adversarial perturbations have long plagued AI systems. Previous work has shown that tiny, imper-
 47 ceptible changes to an image can drastically alter model predictions [34–40]. This issue has driven
 48 extensive research into making AI models more robust, particularly for mission-critical applications.
 49 Yet the impact of adversarial perturbations has received surprisingly little attention in vision neuro-
 50 science. Some work has explored “robustified” encoding models, either through training directly on
 51 neural data [41] to estimate neural robustness or by employing explicit robust pre-training to modify
 52 percepts [42–44]. To our knowledge, no study has directly examined the vulnerability of encoding
 53 models to targeted adversarial perturbations, the relationship between adversarial sensitivity and
 54 predictivity, or the impact of model mapping choices on the model’s adversarial robustness. Under-
 55 standing the bounds of our encoding models is crucial for progress. If imperceptible changes can
 56 distort model predictions, it raises concerns about their reliability in capturing true neural processes
 57 and ability to generalize to unseen data.

58 The central contribution of our work is threefold: (A) we demonstrate that encoding models are
59 susceptible to small-scale adversarial attacks (Figures 1, 3), (B) we show that adversarial sensitivity
60 is a potentially more effective way to differentiate between encoding models than predictive accuracy
61 alone (Figure 3), and (C) we find that the choice of feature-to-brain mapping in encoding model
62 can impact adversarial sensitivity, with sparse mappings producing relatively more robust models of
63 neural activity (Figure 5).

64 2 Methods

65 2.1 fMRI Dataset

66 We used publicly available 7T fMRI data from the Natural Scenes Dataset (NSD) [45] for all analyses
67 in this study. Specifically, we focused on the responses to 515 shared stimuli obtained from fMRI
68 scans of eight subjects in category-selective brain regions. Each subject viewed these images three
69 times over multiple experimental sessions. All analyses were conducted using version 3 of the dataset
70 (betas_fithrf_GLMdenoise_RR), obtained directly from the NSD website. In this work we focused
71 on the category-selective areas: fusiform face area (FFA) [46], extrastriate body area (EBA) [47],
72 parahippocampal place area (PPA) [48], and the visual word form area (VWFA) [49]. To ensure the
73 inclusion of only the most category-selective voxels, we applied a stringent threshold of $tval > 7$ for
74 all analyses. Models were trained to predict the voxel and trial-averaged responses across subjects, as
75 in previous work [9].

76 2.2 Encoding Model

77 Typical ANN-based encoding models consist of two components: embeddings from a specific layer of
78 the artificial neural network (serving as the representational basis) and a trainable mapping function.
79 This mapping is typically done through regularized linear regression, which projects the features into
80 the response subspace of neural activity.

81 Formally, we input each of our training images (see cross-validation schema next) into a represen-
82 tational encoder f and extract the latent feature vector $z_l \in \mathbb{R}^{C_l \times H_l \times W_l}$. These features are then
83 passed through our mapping function $g : \mathbb{R}^{C_l \times H_l \times W_l} \rightarrow \mathbb{R}^n$, where n is the dimensionality of the
84 predicted neural data. To build the encoder model, we freeze f (the weights of the representational
85 encoder) and train the mapping g .

86 **Model architectures:** We considered eight pre-trained artificial neural network architectures that have
87 been previously validated against brain data. These include ResNet-50 [50], VGG16 [51], Inception
88 v3 [52], SqueezeNet v1 [53], AlexNet [54], CORnet-RT [55], DenseNet [56], and MobileNet-v2
89 [57].

90 To investigate whether increasing robustness improves the prediction accuracy of the encoding models,
91 we also used publicly available models that were robustified through adversarial training [58]. These
92 models share the same architecture (ResNet-50) and learning rule but differ in the degree to which
93 they are trained adversarially. More details on the robust models and their training can be found in
94 [59].

95 **Encoding model mapping procedures:** In this study, we experiment with five different mapping
96 functions: ordinary least squares regression (OLS), lasso regression, ridge regression, a two-layer
97 multi-layer perceptron (MLP), and a convolutional neural network (CNN). The first three mapping
98 functions generate direct brain predictions, while the latter two involve learning at least one additional
99 layer of features. These new features may enhance the model’s ability to predict brain responses and
100 could provide more representational robustness. However, the regression methods are computationally
101 faster and do not require extensive hyperparameter tuning for convergence. Our two-layer MLP and
102 CNN both include one hidden layer with 128 units. Note that we used OLS mapping for the first half
103 of the paper because it is the most computationally efficient and does not rely on any assumptions.

104 **Encoding model cross-validation procedure:** We used the 515 shared images across all 8 subjects
105 from the NSD dataset. We trained the model on a randomly chosen set of 400 images and all results
106 in the study are based on predicted responses based on the held-out 115 images.

107 In Section 3.5, we investigate the effect of L_1 readout regularization on the adversarial robustness of
108 the encoding model. We fit each model to the data using only one randomly chosen subject (subj2),

109 testing six different values of the regularization coefficient α (0.0001, 0.001, 0.005, 0.01, 0.05, 0.1).
110 The α value that maximized predictive accuracy for this subject was selected for further analysis. Im-
111 portantly, all model evaluations were conducted using an independent metric (adversarial sensitivity)
112 and across all subjects.

113 2.3 Evaluating encoding model robustness

We evaluated adversarial robustness against the Fast Gradient Sign Method (FGSM) [35]. FGSM attacks are bounded by the L_∞ norm. That is, we find the maximum change δ (bounded by a “perturbation budget” ϵ) predicted to change the response of a given voxel. A successful attack would drastically (and unrealistically) change the predicted response of the encoding model. We quantified the adversarial sensitivity s_i for a given voxel using the method described in [41]. Specifically, we use a sensitivity metric s_i defined as:

$$s_i = \max_{\|\delta\| \leq \epsilon} (g(f(x)) - g(f(x + \delta)))$$

114 There are two things to note about this metric. First, since s_i is a measure of model *sensitivity*, high
115 values on this metric would indicate lower adversarial robustness. We indicate this in several of our
116 plots. The second is that since the metric does not have an upper bound, the results must not be
117 interpreted across regions. While other forms of adversarial attacks exist in the literature, we focus
118 on FGSM for simplicity and consistency.

119 2.4 Encoding Model Discriminability

120 We evaluate the ability of both metrics – adversarial robustness and model predictivity – to dis-
121 criminate encoding models of the brain. For each of the eight models evaluated, we compute
122 the average sensitivity across all subjects and brain regions. We explore whether the spread of the
123 adversarial robustness distribution of the encoding models will be greater than the spread of the
124 model predictivity distribution (i.e., “adversarial robustness” serves as a better discriminative tool).
125 To evaluate this, we test the variance and sparseness of both adversarial sensitivity and predictivity.

126 **Normalized Variance:** Since the scale of “sensitivity” (unbounded) and “predictivity” (bounded -1
127 to 1) are different, we cannot directly compare the variances. Instead, we first divide all accuracy
128 and sensitivity values by their respective maximum value before reporting the variances (hence
129 normalized variance).

130 **Sparseness:** We use the sparseness metric defined in [60, 61]. Specifically, for a distribution of
131 values $P(r)$, sparseness (S) is computed with the following:

$$S = 1 - \frac{E[r]^2}{E[r^2]},$$

132 where $E[\cdot]$ denotes the expectation operator.

133 3 Results

134 Our investigation focuses on category-selective regions—specifically face, body, scene, and word-
135 selective areas (FFA, EBA, PPA, and VWFA, respectively) from the Natural Scenes Dataset (NSD).
136 These regions were chosen because of the extensive work on developing encoding models for them and
137 because they provide the necessary foundational intuition for interpreting changes due to adversarial
138 perturbations (Figure 1C, 1D). We specifically focus on *very small image perturbations* ($\epsilon \leq 3/255$)
139 lower than the resolution limit of the human eye and hence imperceptible to humans. This is because
140 the response of brain voxels to these targeted noise patterns remains currently unknown. Restricting
141 our analysis to small magnitudes ensures that the adversarial sensitivities we detect are real and
142 meaningful.

143 3.1 Several ANN-based encoding models predict voxel responses equally well

144 We first set out to replicate the previous finding that encoding models exhibit similar accuracy in
145 predicting brain responses. To do this, we examined eight pre-trained neural network architectures

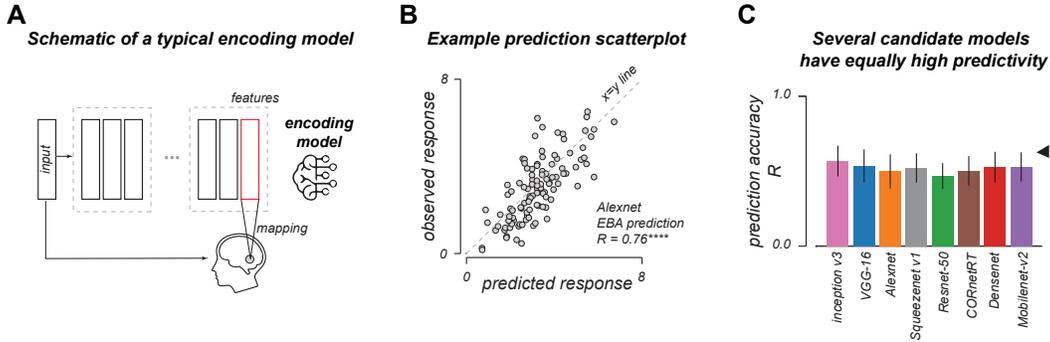


Figure 2: **Several encoding models have equally high predictivity on fMRI data** **A.** Schematic outlining the construction of a typical encoding model. Features from an intermediate model layer (shown in red) are used to build a linear mapping function (indicated as “mapping”) to predict responses in specific brain regions. **B.** Example scatterplot showing predicted (x-axis) versus observed responses for AlexNet in the EBA. The dotted line represents the $x = y$ line, and each dot corresponds to a stimulus that was not used in model training (cross-validated). **C.** Bar plot showing various candidate encoding model architectures (x-axis) and their ability to predict responses to unseen images (y-axis). The black sideways triangle indicates the ceiling performance (median Spearman-Brown corrected split-half correlation across subjects and brain regions). Bars represent the mean response, with error bars showing the SEM across models and brain regions.

146 that have been reported extensively in prior work [14, 10, 9]. For each model, we focused on
 147 features from an intermediate layer, selecting the layer that had previously been shown to achieve the
 148 highest cross-validated accuracy in predicting responses from category-selective regions based on
 149 an independent fMRI dataset [9]. This choice removed experimenter degrees of freedom. Next, we
 150 constructed encoding models by mapping the features from a subset of images to brain responses
 151 using a linear mapping function (see Methods for details on cross-validation procedures). This entire
 152 process is depicted schematically in Figure 2A.

153 Overall, we found that these ANN-based encoding models were highly effective at predicting brain
 154 responses to held-out images (replicating previous findings [10, 9]). This is illustrated for an example
 155 brain region (EBA) in Figure 2B ($R = 0.76$, $P < 0.00001$). Across all regions we considered, the
 156 models were able to predict nearly all of the explainable variance in the observed data. The prediction
 157 accuracy for each model architecture (Figure 2C, bars) was very close to the estimated noise ceiling
 158 (Figure 2C, sideways triangle, derived from corrected split-half correlations). Importantly, all models
 159 appeared to perform similarly well at predicting responses to unseen images. These results replicate
 160 the earlier observation that a wide range of encoding models are approximately equal in their ability
 161 to predict responses in the brain.

162 3.2 All ANN-based encoding models are susceptible to small scale adversarial attacks

163 How susceptible are encoding models to adversarial attacks? To address this, we engineered an
 164 imperceptible noise pattern specifically designed to alter the predicted response for a given brain
 165 region, along with a randomized noise pattern of the same magnitude and statistical properties as
 166 a control. We discovered that even the slightest targeted noise, unseen by the human eye, could
 167 completely derail the encoding model’s predicted response. This is shown for an example encoding
 168 model (VGG16) for the FFA and PPA in Figures 1C and 1D. Initially, the model’s prediction for the
 169 unaltered image from the preferred category (faces for FFA, scenes for PPA) was high. This agrees
 170 with our expectation about images from the preferred category. However, adding a small amount of
 171 targeted noise was enough to push the predicted response well outside the preferred category range
 172 to the extreme end of the observed response spectrum. As a negative control, we used a shuffled
 173 version of the same targeted noise. Importantly, this shuffled noise pattern, despite having the same
 174 summary statistical properties of the noise, did not alter the predicted response to the same extent
 175 ($\Delta = 0.01$).

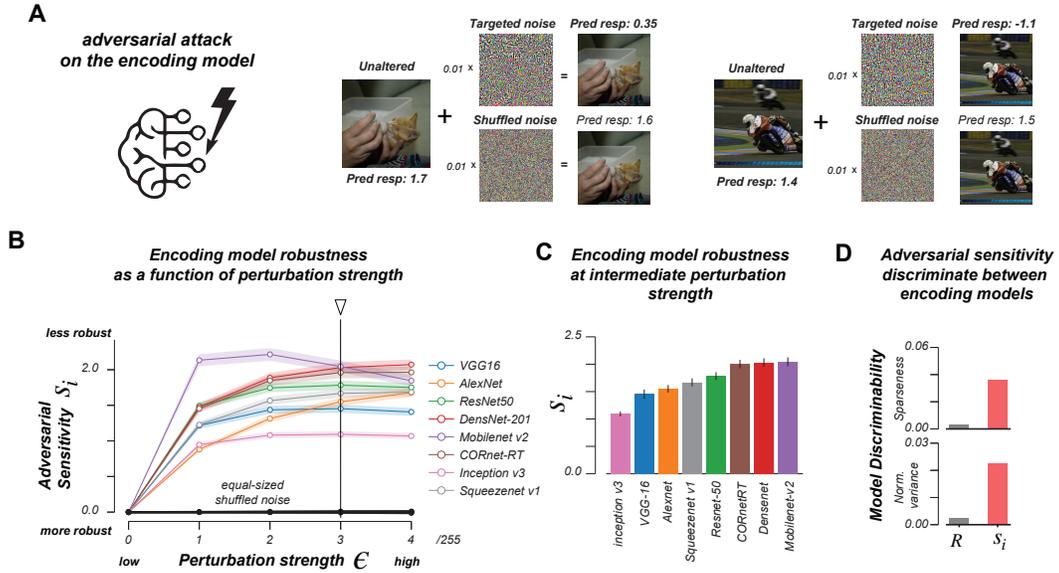


Figure 3: **Adversarial sensitivity effectively discriminates between encoding models.** **A.** Example of adversarial attacks applied directly to the encoding models. Each attack shows the unaltered image (left), the targeted noise and shuffled noise (middle panels top and bottom respectively), and model predictions for these images **B.** The effect of perturbation strength (x-axis) on the model’s adversarial sensitivity (y-axis). Each colored line represents a different candidate encoding model architecture. The dots show the mean sensitivity, and the shaded areas represent the standard error across subjects and brain regions. The black line indicates the negative control using randomized noise. The triangle above marks the perturbation strength used for the subsequent analyses. **C.** Bar plots showing the adversarial sensitivity (y-axis) for all encoding models at a perturbation strength of $3/255$. The models are arranged in the same order as in Figure 2C for direct comparison. **D.** Barplots showing the discriminability between models using adversarial sensitivity and predictivity. Top: Bar plots illustrating model discriminability using predictivity (R) and adversarial sensitivity (s_i). Top: Discriminability based on the sparseness measure (y-axis). Bottom: Discriminability based on a normalized variance measure (y-axis).

176 We quantified the adversarial sensitivity for each model by measuring the change in predicted response
 177 to the adversarially perturbed image. Figure 3B shows these results for all encoding models. As
 178 the strength of the perturbation (ϵ) increased (x-axis), the adversarial sensitivity also increased (as
 179 expected). Note that in this context, higher sensitivity indicates lower adversarial robustness for
 180 the model. These findings demonstrate that all tested ANN models were vulnerable to targeted
 181 adversarial attacks. In fact, for most models, even a small perturbation with $\epsilon = 3/255$ was enough
 182 to significantly alter the predicted response.

183 3.3 Adversarial sensitivity better discriminates between ANN-models than predictivity

184 Next, we evaluated whether adversarial sensitivity could serve as a more effective tool for distin-
 185 guishing between candidate encoding models of the brain. We present these analyses for $\epsilon = 3/255$,
 186 although all subsequent inferences hold across other values as well. The results for adversarial sensi-
 187 tivity across all encoding models at $\epsilon = 3/255$ are displayed in Figure 3C. To facilitate comparison,
 188 the models are arranged in the same order as shown in Figure 2C.

189 To assess the effectiveness of adversarial sensitivity compared to predictivity, we employed two
 190 different measures. First, we measured the sparseness [61] of the adversarial sensitivity and predictiv-
 191 ity metrics across models. Sparseness was chosen since it is a scale invariant measure and can be
 192 used to directly compare between predictivity and adversarial sensitivity (see Methods for details).
 193 Figure 3D (top) shows that model sparseness was significantly higher for adversarial sensitivity
 194 than for predictivity, indicating better discriminability across models. A problem with sparseness
 195 however is that it is highly sensitive to outliers. To allay this concern, we adopted a second, more

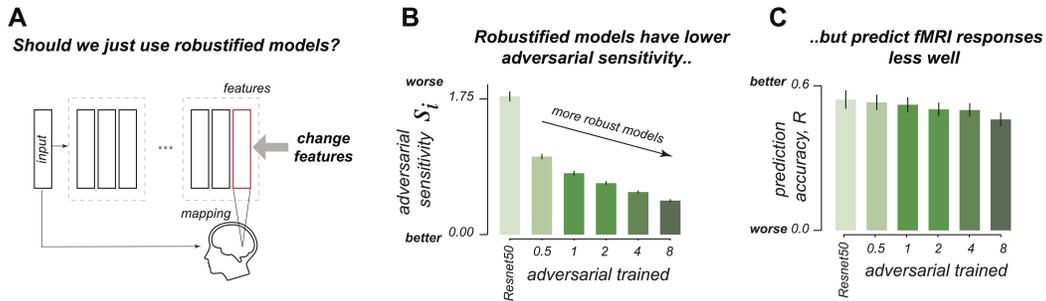


Figure 4: **Robust training reduces the predictive accuracy of fMRI encoding models.** **A.** Schematic illustration of the analysis. In this step, we replace the original features (shown in red) with features that have been robustified through adversarial training. **B.** Bar plots showing that increasing the level of adversarial training (x-axis) improves the adversarial sensitivity of the encoding models (y-axis). **C.** Bar plots showing that increasing the level of adversarial training (x-axis) reduces the predictive accuracy of encoding models on fMRI data (y-axis).

196 intuitive variance measure (normalized to match the scale between sensitivity and predictivity). As
 197 shown in Figure 3D (bottom), the normalized variance was also higher for adversarial sensitivity
 198 compared to predictivity. Together, these measures present a consistent picture: adversarial sensitivity
 199 distinguishes between encoding models more effectively than predictivity alone.

200 3.4 Increasing model robustness via adversarial training does not improve model predictivity

201 So far, we have demonstrated two key findings: 1) commonly used encoding models are sensitive
 202 to imperceptible adversarial noise, and 2) adversarial sensitivity can serve as a tool to distinguish
 203 between predictive models. How can we build better, more robust encoding models? The natural
 204 thing to try is to simply replace the current model architecture with a more robust one. In this section,
 205 we explored what happens when we use robustified models. To test this question, we fixed the model
 206 architecture (ResNet50) and parametrically varied the strength of adversarial training using publicly
 207 available robustified models [58]. This strategy is illustrated schematically in Figure 4A.

208 As expected, we found that robust models were indeed less vulnerable to added adversarial noise.
 209 Figure 4B shows how adversarial sensitivity decreases as the strength of adversarial training increases.
 210 Are robustified models effective at predicting fMRI responses? Here, we observed a trade-off: as
 211 the models became more robust, their ability to predict fMRI responses declined. This reduction
 212 was quite significant and is shown across all models and regions. These results suggest that while
 213 adversarial training does improve robustness, it may do so at the cost of reduced predictivity for
 214 brain data.

215 3.5 Sparse mappings tend to improve adversarial robustness of encoding models without 216 sacrificing model predictivity

217 A less well-understood aspect of encoding models is the effect of the specific choice of mapping
 218 between model features and neural responses. We wondered if certain mapping functions could
 219 improve an encoding model's sensitivity to targeted noise. There are many potential linear and
 220 non-linear mapping functions to explore. To constrain our choices, we first evaluated five different
 221 mapping methods: ordinary least squares (no regularization), Lasso (L_1) regression (sparse), Ridge
 222 (L_2) regression, a two-layer multi-layer perceptron (MLP), and a convolutional neural network.
 223 We chose two candidate encoding models (VGG16 and ResNet50) for this initial exploration of
 224 mapping methods. An issue with these is that many of these methods involve choosing appropriate
 225 hyperparameters. Hyperparameters were selected based on prediction accuracy (see Methods), but we
 226 focus our attention on an independent metric: adversarial sensitivity. The results are presented in
 227 Table 1. Across both models, we found evidence of a significant boost in adversarial sensitivity when
 228 using a sparse mapping.

229 Would this observation generalize to other models? To explore this, we compared the sensitivity of all
 230 eight models using L_1 (sparse) and OLS (no regularization) mapping-based encoding models across

Adversarial sensitivity for different model-to-brain mapping functions

Model	OLS	L1 (Lasso)	L2 (Ridge)	2-layer MLP	CNN
VGG16	1.453	.891	1.453	1.358	1.734
ResNet50	1.782	.821	1.782	1.286	1.051

Table 1: Effect of readout functions on adversarial sensitivity. L_1 regularization on the readout performed best. The weight of the regularization term, α , was chosen as the value which maximized predictive accuracy from a set of candidate values; see Methods.

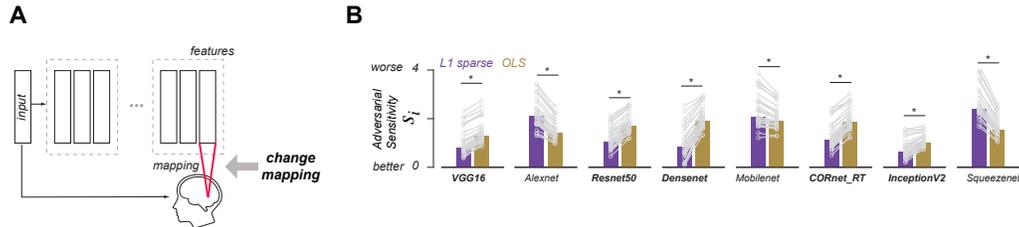


Figure 5: Sparse model-to-brain mappings tend to lower adversarial sensitivity of encoding models . A. Schematic illustration of the analysis. Here, we evaluate the model-to-brain mappings (highlighted in red) **B.** Bar plots showing adversarial sensitivity (y-axis) for all models tested. The dots and connected lines represent an encoding model for a specific subject and brain region. * indicates statistical significance (paired t-test, $P < 0.001$) between OLS and sparse mappings. Models in bold indicate improved adversarial sensitivity for sparse (L_1) mappings compared to OLS-based mappings.

231 all architectures. Note that the hyperparameters here were determined based on prediction accuracy by
 232 from one subject, and the results were independently evaluated on adversarial sensitivity from all
 233 subjects (see Methods). This preliminary analysis revealed an interesting trend: sparse mappings
 234 produced significantly more robust models in 5 out of 8 model architectures. While this suggests that
 235 sparse mappings may enhance adversarial robustness, it is important to emphasize that these results
 236 are still preliminary and additional testing is needed to confirm whether this pattern holds across a
 237 larger sample, different model types, and independent analysis methods. Nonetheless, these early
 238 findings hint at the potential of sparse mappings to provide a meaningful boost in robustness.

239 4 Discussion

240 In this study, we investigated how susceptible commonly used ANN-based vision encoding models
 241 were to small-scale adversarial perturbations. We found that all high-performing models were
 242 vulnerable to imperceptible, small-scale adversarial noise (Figure 3). We also demonstrated that
 243 adversarial sensitivity, more effectively than prediction accuracy, could be used to differentiate
 244 between models (Figure 3). However, increasing model robustness through adversarial training came
 245 at the expense of reducing their ability to predict fMRI responses (Figure 4). Finally, we found early
 246 evidence that a simple sparse mapping approach on the mapping function could significantly improve
 247 adversarial robustness (Figure 5). These findings reveal key limitations of current encoding models
 248 and suggest new strategies for enhancing their performance.

249 Our adversarial attacks had two key features. First, the perturbations were deliberately kept small to
 250 focus on imperceptible changes. Our pilot analyses, based on an 8-degree viewing angle, suggest the
 251 detection threshold for adversarial images to be around $\epsilon = 8/255$. While a formal estimate on a
 252 larger sample is underway, we assumed that small perturbations, as those used in this study, would
 253 not alter brain voxel responses (though see [41]). This allowed us to test the model’s vulnerability
 254 in a regime where the visual system should remain stable, highlighting its susceptibility to subtle
 255 adversarial noise. However, these assumptions require formal testing in future work. The second
 256 key feature is that our method targeted the encoding models directly (instead of the model features).
 257 This approach enabled us to assess vulnerabilities in the model’s representational mappings to brain
 258 activity, not just the image embeddings. While previous studies have examined the relationship
 259 between model robustness and neural predictions in monkeys [44], or the link between spatial features

260 and neural representations [62, 63], our work extends these findings by exploring how adversarial
261 perturbations *directly* affect model representations most predictive of human fMRI brain responses.

262 One interpretation of our results is that current high-performing, predictive encoding models are
263 fundamentally flawed given how drastically they fail when exposed to targeted adversarial noise.
264 While this is true, our aim is not merely to highlight these vulnerabilities. It is not entirely unexpected
265 that these models are susceptible to adversarial perturbations, given what we know about neural
266 networks in general. However, we propose leveraging adversarial sensitivity as a tool to guide the
267 development of more accurate and resilient models. In fact, we find that adversarial sensitivity
268 provides an additional layer of insight into model performance, helping to distinguish between highly
269 predictive encoding models.

270 By analyzing how different models respond to adversarial perturbations, we start to uncover their
271 limitations and use new insights into the development of more robust brain models. To this end, we
272 tested two strategies. While adversarial training is widely used in the AI community to enhance
273 model resilience, we found that it came at a significant cost to model predictivity (see also [44]). As
274 models became more robust, their ability to accurately predict brain responses declined substantially.
275 This trade-off highlights a compromise that must be carefully considered when developing models
276 for neuroscience applications. In contrast, we found that a relatively simple sparse mapping between
277 model features and brain representations was enough to significantly reduce the adversarial sensitivity
278 of most encoding models, usually outperforming more complex non-linear mapping methods. We
279 hope to explore these differences further in future work.

280 Taken together, our results expose the critical vulnerabilities of ANN-based predictive encoding mod-
281 els to adversarial perturbations, highlight adversarial sensitivity as a powerful tool for differentiating
282 between models, and suggest a promising path for enhancing model robustness. As we continue
283 our search for brain-like models, striking the right balance between robustness and predictivity will
284 be crucial. Our work provides a foundation for tracking this balance, offers new model evaluations,
285 and offers prescriptions to guide the development of more accurate and resilient models that can be
286 applied to study human cognition even beyond vision.

287 **References**

- 288 [1] Martin Schrimpf, Jonas Kubilius, Michael J Lee, N Apurva Ratan Murty, Robert Ajemian, and James J
289 DiCarlo. Integrative benchmarking to advance neurally mechanistic models of human intelligence. *Neuron*,
290 2020.
- 291 [2] Katherine R Storrs, Tim C Kietzmann, Alexander Walther, Johannes Mehrer, and Nikolaus Kriegeskorte.
292 Diverse deep neural networks all predict human inferior temporal cortex well, after training and fitting.
293 *Journal of cognitive neuroscience*, 33(10):2044–2064, 2021.
- 294 [3] Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C Frank, James J DiCarlo, and
295 Daniel LK Yamins. Unsupervised neural network models of the ventral visual stream. *Proceedings of the*
296 *National Academy of Sciences*, 118(3):e2014196118, 2021.
- 297 [4] Nikolaus Kriegeskorte. Deep neural networks: a new framework for modeling biological vision and brain
298 information processing. *Annual review of vision science*, 1(1):417–446, 2015.
- 299 [5] Adrien Doerig, Rowan P Sommers, Katja Seeliger, Blake Richards, Jenann Ismael, Grace W Lindsay, Kon-
300 rad P Kording, Talia Konkle, Marcel AJ Van Gerven, Nikolaus Kriegeskorte, et al. The neuroconnectionist
301 research programme. *Nature Reviews Neuroscience*, 24(7):431–450, 2023.
- 302 [6] Blake A Richards, Timothy P Lillicrap, Philippe Beaudoin, Yoshua Bengio, Rafal Bogacz, Amelia
303 Christensen, Claudia Clopath, Rui Ponte Costa, Archy de Berker, Surya Ganguli, et al. A deep learning
304 framework for neuroscience. *Nature neuroscience*, 22(11):1761–1770, 2019.
- 305 [7] Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo.
306 Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings*
307 *of the national academy of sciences*, 111(23):8619–8624, 2014.
- 308 [8] Umut Güçlü and Marcel AJ Van Gerven. Deep neural networks reveal a gradient in the complexity of
309 neural representations across the ventral stream. *Journal of Neuroscience*, 35(27):10005–10014, 2015.
- 310 [9] N. Apurva Ratan Murty, Pouya Bashivan, Alex Abate, James J. DiCarlo, and Nancy Kanwisher. Com-
311 putational models of category-selective brain regions enable high-throughput tests of selectivity. *Nature*
312 *Communications*, 12, Sep 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-25409-6.
- 313 [10] Colin Conwell, Jacob S. Prince, Kendrick N. Kay, George A. Alvarez, and Talia Konkle. What can
314 1.8 billion regressions tell us about the pressures shaping high-level visual representation in brains and
315 machines? *bioRxiv*, 2023. doi: 10.1101/2022.03.28.485868.
- 316 [11] Charles Y. Zheng, Francisco Pereira, Chris I. Baker, and Martin N. Hebart. Revealing interpretable object
317 representations from human behavior. In *International Conference on Learning Representations*, 2019.
- 318 [12] Rishi Rajalingham, Elias B Issa, Pouya Bashivan, Kohitij Kar, Kailyn Schmidt, and James J DiCarlo.
319 Large-scale, high-resolution comparison of the core visual object recognition behavior of humans, monkeys,
320 and state-of-the-art deep artificial neural networks. *Journal of Neuroscience*, 38(33):7255–7269, 2018.
- 321 [13] Katharina Dobs, Joanne Yuan, Julio Martinez, and Nancy Kanwisher. Behavioral signatures of face
322 perception emerge in deep neural networks optimized for face recognition. *Proceedings of the National*
323 *Academy of Sciences*, 120(32):e2220642120, 2023.
- 324 [14] Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J. Majaj, Rishi Rajalingham, Elias B. Issa, Kohitij Kar,
325 Pouya Bashivan, Jonathan Prescott-Roy, Franziska Geiger, Kailyn Schmidt, Daniel L. K. Yamins, and
326 James J. DiCarlo. Brain-score: Which artificial neural network for object recognition is most brain-like?
327 *bioRxiv preprint*, 2018.
- 328 [15] Drew Linsley, Ivan F. Rodriguez Rodriguez, Thomas Fel, Michael Arcaro, Saloni Sharma, Margaret S.
329 Livingstone, and Thomas Serre. Performance-optimized deep neural networks are evolving into worse
330 models of inferotemporal visual cortex. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko,
331 Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual*
332 *Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA,*
333 *December 10 - 16, 2023*, 2023.
- 334 [16] Jeffrey S Bowers, Gaurav Malhotra, Marin Dujmović, Milton Llera Montero, Christian Tsvetkov, Valerio
335 Biscione, Guillermo Puebla, Federico Adolffi, John E Hummel, Rachel F Heaton, et al. Deep problems
336 with neural network models of human vision. *Behavioral and Brain Sciences*, 46:e385, 2023.

- 337 [17] Jonas Kubilius, Martin Schrimpf, Ha Hong, Najib J. Majaj, Rishi Rajalingham, Elias B. Issa, Kohitij Kar,
338 Pouya Bashivan, Jonathan Prescott-Roy, Kailyn Schmidt, Aran Nayebi, Daniel Bear, Daniel L. K. Yamins,
339 and James J. DiCarlo. Brain-Like Object Recognition with High-Performing Shallow Recurrent ANNs.
340 In H. Wallach, H. Larochelle, A. Beygelzimer, F. D'Alché-Buc, E. Fox, and R. Garnett, editors, *Neural*
341 *Information Processing Systems (NeurIPS)*, 2019.
- 342 [18] Tim C. Kietzmann, Courtney J. Sporer, Lynn K. A. Sörensen, Radoslaw M. Cichy, Olaf Hauk, and
343 Nikolaus Kriegeskorte. Recurrence is required to capture the representational dynamics of the human
344 visual system. *Proceedings of the National Academy of Sciences*, 116(43):21854–21863, 2019. doi:
345 10.1073/pnas.1905544116.
- 346 [19] Ruben S van Bergen and Nikolaus Kriegeskorte. Going in circles is the way forward: the role of recurrence
347 in visual inference. *Current Opinion in Neurobiology*, 65:176–193, 2020.
- 348 [20] Aran Nayebi, Daniel Bear, Jonas Kubilius, Kohitij Kar, Surya Ganguli, David Sussillo, James J DiCarlo,
349 and Daniel L Yamins. Task-driven convolutional recurrent models of the visual system. *Advances in neural*
350 *information processing systems*, 31, 2018.
- 351 [21] Thomas Fel, Ivan F. Rodriguez Rodriguez, Drew Linsley, and Thomas Serre. Harmonizing the object
352 recognition strategies of deep neural networks with humans. In Sanmi Koyejo, S. Mohamed, A. Agarwal,
353 Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35:*
354 *Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA,*
355 *USA, November 28 - December 9, 2022*, 2022.
- 356 [22] Joel Dapello, Kohitij Kar, Martin Schrimpf, Robert Geary, Michael Ferguson, David D Cox, and James J
357 DiCarlo. Aligning model and macaque inferior temporal cortex representations improves model-to-human
358 behavioral alignment and adversarial robustness. *bioRxiv*, pages 2022–07, 2022.
- 359 [23] Meenakshi Khosla and Leila Wehbe. High-level visual areas act like domain-general filters with strong
360 selectivity and functional specialization. *bioRxiv*, pages 2022–03, 2022.
- 361 [24] Yena Han, Tomaso A. Poggio, and Brian Cheung. System identification of neural systems: If we got it
362 right, would we know? In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan
363 Sabato, and Jonathan Scarlett, editors, *International Conference on Machine Learning, ICML 2023, 23-29*
364 *July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages
365 12430–12444. PMLR, 2023.
- 366 [25] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey E. Hinton. Similarity of neural network
367 representations revisited. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the*
368 *36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California,*
369 *USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 3519–3529. PMLR, 2019.
- 370 [26] Ilya Sucholutsky, Lukas Muttenthaler, Adrian Weller, Andi Peng, Andreea Bobu, Been Kim, Bradley C.
371 Love, Erin Grant, Jascha Achterberg, Joshua B. Tenenbaum, Katherine M. Collins, Katherine L. Hermann,
372 Kerem Oktar, Klaus Greff, Martin N. Hebart, Nori Jacoby, Qiuyi Zhang, Raja Marjeh, Robert Geirhos,
373 Sherol Chen, Simon Kornblith, Sunayana Rane, Talia Konkle, Thomas P. O’Connell, Thomas Unterthiner,
374 Andrew K. Lampinen, Klaus-Robert Müller, Mariya Toneva, and Thomas L. Griffiths. Getting aligned on
375 representational alignment. *CoRR*, abs/2310.13018, 2023.
- 376 [27] Joel Dapello, Kohitij Kar, Martin Schrimpf, Robert Baldwin Geary, Michael Ferguson, David Daniel
377 Cox, and James J. DiCarlo. Aligning model and macaque inferior temporal cortex representations
378 improves model-to-human behavioral alignment and adversarial robustness. In *The Eleventh International*
379 *Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net,
380 2023.
- 381 [28] Meenakshi Khosla and Alex H. Williams. Soft matching distance: A metric on neural representations
382 that captures single-neuron tuning. In Marco Fumero, Emanuele Rodolà, Clémentine Dominé, Francesco
383 Locatello, Karolina Dziugaite, and Mathilde Caron, editors, *Proceedings of UniReps: the First Workshop*
384 *on Unifying Representations in Neural Models, 15 December 2023, Ernest N. Morial Convention Center,*
385 *New Orleans, USA*, volume 243 of *Proceedings of Machine Learning Research*, pages 326–341. PMLR,
386 2023.
- 387 [29] Meenakshi Khosla and Leila Wehbe. High-level visual areas act like domain-general filters with strong
388 selectivity and functional specialization. *bioRxiv*, 2022. doi: 10.1101/2022.03.16.484578.
- 389 [30] Rishi Rajalingham, Elias B Issa, Pouya Bashivan, Kohitij Kar, Kailyn Schmidt, and James J DiCarlo.
390 Large-scale, high-resolution comparison of the core visual object recognition behavior of humans, monkeys,
391 and state-of-the-art deep artificial neural networks. *J. Neurosci.*, 38(33):7255–7269, August 2018.

- 392 [31] Pouya Bashivan, Kohitij Kar, and James J DiCarlo. Neural population control via deep image synthesis.
393 *Science*, 364(6439):eaav9436, 2019.
- 394 [32] Subha Nawer Pushpita, Elizabeth Mieczkowski, Bradley Duchaine, and N. Apurva Ratan Murty. Intensive
395 fmri scanning and computational models can provide insight into the neural basis of developmental
396 prosopagnosia. *Journal of Vision*, 23(9):5838, 2023. doi: <https://doi.org/10.1167/jov.23.9.5838>.
- 397 [33] Carlos R. Ponce, Will Xiao, Peter F. Schade, Till S. Hartmann, Gabriel Kreiman, and Margaret S.
398 Livingstone. Evolving images for visual neurons using a deep generative network reveals coding
399 principles and neuronal preferences. *Cell*, 177(4):999–1009.e10, 2019. ISSN 0092-8674. doi:
400 <https://doi.org/10.1016/j.cell.2019.04.005>.
- 401 [34] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and
402 Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun, editors, *2nd*
403 *International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014,*
404 *Conference Track Proceedings*, 2014.
- 405 [35] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples.
406 In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations,*
407 *ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- 408 [36] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *5th*
409 *International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017,*
410 *Workshop Track Proceedings*, 2017.
- 411 [37] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of
412 diverse parameter-free attacks. In *ICML*, 2020.
- 413 [38] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural
414 networks. *IEEE Trans. Evol. Comput.*, 23(5):828–841, 2019.
- 415 [39] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adver-
416 sarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
417 pages 1765–1773, 2017.
- 418 [40] Yatie Xiao, Chi-Man Pun, and Kongyang Chen. Towards evaluating the robustness of deep neural semantic
419 segmentation networks with feature-guided method. *Knowl. Based Syst.*, 281:111063, 2023.
- 420 [41] Chong Guo, Michael J. Lee, Guillaume Leclerc, Joel Dapello, Yug Rao, Aleksander Madry, and James J.
421 DiCarlo. Adversarially trained neural representations may already be as robust as corresponding biological
422 neural representations. In *39th International Conference on Machine Learning, ICML 2022, Baltimore,*
423 *MD, USA, 2015, Conference Track Proceedings*, 2022.
- 424 [42] Guy Gaziv, Michael J. Lee, and James J. DiCarlo. Strong and precise modulation of human percepts via
425 robustified ans. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey
426 Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*
427 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023,*
428 2023.
- 429 [43] Hojin Jang and Frank Tong. Improved modeling of human vision by incorporating robustness to blur in
430 convolutional neural networks. *Nature Communications*, 15(1):1989, 2024.
- 431 [44] Yifei Ren and Pouya Bashivan. How well do models of visual cortex generalize to out of distribution
432 samples? *PLOS Computational Biology*, 20(5):e1011145, 2024.
- 433 [45] Emily J. Allen, Ghislain St-Yves, Yihan Wu, Jesse L. Breedlove, Logan T. Dowdle, Brad Caron, Franco
434 Pestilli, Ian Charest, J. Benjamin Hutchinson, Thomas Naselaris, and Kendrick Kay. A massive 7t
435 fmri dataset to bridge cognitive and computational neuroscience. *Nature Neuroscience*, 2022. doi:
436 [10.1038/s41593-021-00962-x](https://doi.org/10.1038/s41593-021-00962-x).
- 437 [46] Nancy Kanwisher, Josh McDermott, and Marvin M. Chun. The fusiform face area: a module in human
438 extrastriate cortex specialized for face perception. *J. Neurosci.*, 17(11):4302–4311, June 1997.
- 439 [47] P E Downing, Yuhong Jiang, M Shuman, and N Kanwisher. A cortical area selective for visual processing
440 of the human body. *Science (New York, N.Y.)*, 293(5539):2470–3, September 2001. ISSN 0036-8075.
- 441 [48] Russell A. Epstein and Nancy G. Kanwisher. A cortical representation of the local visual environment.
442 *Nature*, 392:598–601, 1998.

- 443 [49] Laurent Cohen, Stanislas Dehaene, Lionel Naccache, Stéphane Lehericy, Ghislaine Dehaene-Lambertz,
444 Marie-Anne Hénaff, and François Michel. The visual word form area: Spatial and temporal characterization
445 of an initial stage of reading in normal subjects and posterior split-brain patients. *Brain*, 123(2):291–307,
446 02 2000. ISSN 0006-8950.
- 447 [50] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
448 In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA,
449 June 27-30, 2016*, pages 770–778. IEEE Computer Society, 2016. doi: 10.1109/CVPR.2016.90.
- 450 [51] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recogni-
451 tion. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representa-
452 tions, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- 453 [52] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking
454 the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern
455 Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2818–2826. IEEE Computer
456 Society, 2016.
- 457 [53] Forrest N. Iandola, Matthew W. Moskewicz, Khalid Ashraf, Song Han, William J. Dally, and Kurt
458 Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size. *CoRR*,
459 abs/1602.07360, 2016.
- 460 [54] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional
461 neural networks. In Peter L. Bartlett, Fernando C. N. Pereira, Christopher J. C. Burges, Léon Bottou,
462 and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25: 26th Annual
463 Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6,
464 2012, Lake Tahoe, Nevada, United States*, pages 1106–1114, 2012.
- 465 [55] Jonas Kubilius, Martin Schrimpf, Aran Nayebi, Daniel Bear, Daniel LK Yamins, and James J DiCarlo.
466 Cornet: Modeling the neural mechanisms of core object recognition. *bioRxiv*, 2018.
- 467 [56] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected
468 convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR
469 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2261–2269. IEEE Computer Society, 2017.
- 470 [57] Mark Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo-
471 bilenetv2: Inverted residuals and linear bottlenecks. In *2018 IEEE Conference on Computer Vision and
472 Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 4510–4520. Computer
473 Vision Foundation / IEEE Computer Society, 2018.
- 474 [58] Logan Engstrom, Andrew Ilyas, Hadi Salman, Shibani Santurkar, and Dimitris Tsipras. Robustness (python
475 library), 2019. URL <https://github.com/MadryLab/robustness>.
- 476 [59] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry.
477 Adversarial examples are not bugs, they are features. *Advances in neural information processing systems*,
478 32, 2019.
- 479 [60] Ben D B Willmore, James A Mazer, and Jack L Gallant. Sparse coding in striate and extrastriate visual
480 cortex. *J. Neurophysiol.*, 105(6):2907–2919, June 2011.
- 481 [61] William E. Vinje and Jack L. Gallant. Sparse coding and decorrelation in primary visual cortex during
482 natural vision. *Science*, 287(5456):1273–1276, 2000. doi: 10.1126/science.287.5456.1273.
- 483 [62] Zhe Li, Josue Ortega Caro, Evgenia Rusak, Wieland Brendel, Matthias Bethge, Fabio Anselmi, Ankit B
484 Patel, Andreas S Tolias, and Xaq Pitkow. Robust deep learning object recognition models rely on low
485 frequency information in natural images. *PLOS Computational Biology*, 19(3):e1010932, 2023.
- 486 [63] Ajay Subramanian, Elena Sizikova, Najib J. Majaj, and Denis G. Pelli. Spatial-frequency channels, shape
487 bias, and adversarial robustness. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz
488 Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual
489 Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA,
490 December 10 - 16, 2023*, 2023.