

---

# Overinterpretation reveals image classification model pathologies

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Image classifiers are typically scored on their test set accuracy, but high accuracy  
2 can mask a subtle type of model failure. We find that high scoring convolutional  
3 neural networks (CNNs) on popular benchmarks exhibit troubling pathologies  
4 that allow them to display high accuracy even in the absence of semantically  
5 salient features. When a model provides a high-confidence decision without salient  
6 supporting input features, we say the classifier has overinterpreted its input, finding  
7 too much class-evidence in patterns that appear nonsensical to humans. Here, we  
8 demonstrate that neural networks trained on CIFAR-10 and ImageNet suffer from  
9 overinterpretation, and we find models on CIFAR-10 make confident predictions  
10 even when 95% of input images are masked and humans cannot discern salient  
11 features in the remaining pixel-subsets. We introduce Batched Gradient SIS, a  
12 new method for discovering sufficient input subsets for complex datasets, and use  
13 this method to show the sufficiency of border pixels in ImageNet for training and  
14 testing. Although these patterns portend potential model fragility in real-world  
15 deployment, they are in fact valid statistical patterns of the benchmark that alone  
16 suffice to attain high test accuracy. Unlike adversarial examples, overinterpretation  
17 relies upon unmodified image pixels. We find ensembling and input dropout can  
18 each help mitigate overinterpretation.

## 19 1 Introduction

20 Well-founded decisions by machine learning (ML) systems are critical for high-stakes applications  
21 such as autonomous vehicles and medical diagnosis. Pathologies in models and their respective  
22 training datasets can result in unintended behavior during deployment if the systems are confronted  
23 with novel situations. For example, a medical image classifier for cancer detection attained high  
24 accuracy in benchmark test data, but was found to base decisions upon presence of rulers in an image  
25 (present when dermatologists already suspected cancer) [1]. We define model *overinterpretation* to  
26 occur when a classifier finds strong class-evidence in regions of an image that contain no semantically  
27 salient features. Overinterpretation is related to overfitting, but overfitting can be diagnosed via  
28 reduced test accuracy. Overinterpretation can stem from true statistical signals in the underlying  
29 dataset distribution that happen to arise from particular properties of the data source (e.g., derma-  
30 tologists' rulers). Thus, overinterpretation can be harder to diagnose as it admits decisions that are  
31 made by statistically valid criteria, and models that use such criteria can excel at benchmarks. We  
32 demonstrate overinterpretation occurs with unmodified subsets of the original images. In contrast  
33 to *adversarial examples* that modify images with extra information, overinterpretation is based on  
34 real patterns already present in the training data that also generalize to the test distribution. Hidden  
35 statistical signals of benchmark datasets can result in models that overinterpret or do not generalize  
36 to new data from a different distribution. Computer vision (CV) research relies on datasets like

37 CIFAR-10 [2] and ImageNet [3] to provide standardized performance benchmarks. Here, we analyze  
38 the overinterpretation of popular CNN architectures on these benchmarks to characterize pathologies.

39 Revealing overinterpretation requires a systematic way to identify which features are used by a model  
40 to reach its decision. Feature attribution is addressed by a large number of interpretability methods,  
41 although they propose differing explanations for the decisions of a model. One natural explanation  
42 for image classification lies in the set of pixels that is sufficient for the model to make a confident  
43 prediction, even in the absence of information about the rest of the image. In the example of the  
44 medical image classifier for cancer detection, one might identify the pathological behavior by finding  
45 pixels depicting the ruler alone suffice for the model to confidently output the same classifications.  
46 This idea of Sufficient Input Subsets (SIS) has been proposed to help humans interpret the decisions  
47 of black-box models [4]. An SIS subset is a minimal subset of features (e.g., pixels) that suffices to  
48 yield a class probability above a certain threshold with all other features masked.

49 We demonstrate that classifiers trained on CIFAR-10 and ImageNet can base their decisions on  
50 SIS subsets that contain few pixels and lack human understandable semantic content. Nevertheless,  
51 these SIS subsets contain statistical signals that generalize across the benchmark data distribution,  
52 and we are able to train classifiers on CIFAR-10 images missing 95% of their pixels and ImageNet  
53 images missing 90% of their pixels with minimal loss of test accuracy. Thus, these benchmarks  
54 contain inherent statistical shortcuts that classifiers optimized for accuracy can learn to exploit,  
55 instead of learning more complex *semantic* relationships between the image pixels and the assigned  
56 class label. While recent work suggests adversarially robust models base their predictions on more  
57 semantically meaningful features [5], we find these models suffer from overinterpretation as well.  
58 As we subsequently show, overinterpretation is not only a conceptual issue, but can actually harm  
59 overall classifier performance in practice. We find model ensembling and input dropout partially  
60 mitigate overinterpretation, increasing the semantic content of the resulting SIS subsets. However,  
61 this mitigation is not a substitute for better training data, and we find that overinterpretation is a  
62 statistical property of common benchmarks. Intriguingly, the number of pixels in the SIS rationale  
63 behind a particular classification is often indicative of whether the image is correctly classified.

64 It may seem unnatural to use an interpretability method that produces feature attributions that look  
65 uninterpretable. However, we do not want to bias extracted rationales towards human visual priors  
66 when analyzing a model’s pathologies, but rather faithfully report the features used by a model. To  
67 our knowledge, this is the first analysis showing one can extract nonsensical features from CIFAR-10  
68 and ImageNet that intuitively should be insufficient or irrelevant for a confident prediction, yet are  
69 alone sufficient to train classifiers with minimal loss of performance. Our contributions include:

- 70 • We discover the pathology of overinterpretation and find it is a common failure mode of ML  
71 models, which latch onto non-salient but statistically valid signals in datasets (Section 4.1).
- 72 • We introduce Batched Gradient SIS, a new masking algorithm to scale SIS to high-  
73 dimensional inputs and apply it to characterize overinterpretation on ImageNet (Section 3.2).
- 74 • We provide a pipeline for detecting overinterpretation by masking over 90% of each image,  
75 demonstrating minimal loss of test accuracy, and establish lack of saliency in these patterns  
76 through human accuracy evaluations (Sections 3.3, 4.2, 4.3).
- 77 • We show misclassifications often rely on smaller and more spurious feature subsets suggest-  
78 ing overinterpretation is a serious practical issue (Section 4.4).
- 79 • We identify two strategies for mitigating overinterpretation (Section 4.5). We demonstrate  
80 that overinterpretation is caused by spurious statistical signals in training data, and thus  
81 training data must be carefully curated to eliminate overinterpretation artifacts.

## 82 2 Related Work

83 While existing work has demonstrated numerous distinct flaws in deep image classifiers our paper  
84 demonstrates a new distinct flaw, overinterpretation, previously undocumented in the literature. There  
85 has been substantial research on understanding dataset bias in CV [6, 7] and the fragility of image  
86 classifiers deployed outside benchmark settings. We extend previous work on sufficient input subsets  
87 (SIS) [4] with the Batched Gradient SIS method, and use this method to show that ImageNet sufficient  
88 input subset pixels for training and testing are typically found at image borders. We comprehensively  
89 contrast overinterpretation against known flaws below.

- 90 • Image classifiers have been shown to be fragile when objects from one image are transplanted  
91 in another image [8], and can be biased by object context [9, 10]. In contrast, overinterpretation  
92 differs because we demonstrate that highly sparse, unmodified subsets of pixels in images suffice  
93 for image classifiers to make the same predictions as on the full images.
- 94 • Lapuschkin et al. [11] demonstrate that DNNs can learn to rely on spurious signals in datasets,  
95 including source tags and artificial padding, but which are still human-interpretable. In contrast, the  
96 patterns we identify are minimal collections of pixels in images that are semantically meaningless  
97 to humans (they do not comprise human-interpretable parts of images). We demonstrate such  
98 patterns generalize to the test distribution suggesting they arise from degenerate signals in popular  
99 benchmarks, and thus models trained on these datasets may fail to generalize to real-world data.
- 100 • CNNs in particular have been conjectured to pick up on localized features like texture instead  
101 of more global features like object shape [12, 13]. Brendel and Bethge [14] show CNNs trained  
102 on natural ImageNet images may rely on local features and, unlike humans, are able to classify  
103 texturized images, suggesting ImageNet alone is insufficient to force DNNs to rely on more causal  
104 representations. Our work demonstrates another source of degeneracy of popular image datasets,  
105 where sparse, unmodified subsets of training images that are meaningless to humans can enable a  
106 model to generalize to test data. We provide one explanation for why ImageNet-trained models  
107 may struggle to generalize to out-of-distribution data.
- 108 • Geirhos et al. [15] discover that DNNs trained on distorted images fail to generalize as well as  
109 human observers when trained under image distortions. In contrast, overinterpretation reveals a  
110 different failure mode of DNNs, whereby models latch onto spurious but statistically valid sets of  
111 features in undistorted images. This phenomenon can limit the ability of a DNN to generalize to  
112 real-world data even when trained on natural images.
- 113 • Other work has shown deep image classifiers can make confident predictions on nonsensical  
114 patterns [16], and the susceptibility of DNNs to adversarial examples or synthetic images has been  
115 widely studied [5, 17-19]. However, these adversarial examples synthesize artificial images or  
116 modify real images with auxiliary information. In contrast, we demonstrate overinterpretation of  
117 unmodified subsets of actual training images, indicating the patterns are already present in the  
118 original dataset. We further demonstrate that such signals in training data actually generalize to the  
119 test distribution and that adversarially robust models also suffer from overinterpretation.
- 120 • Hooker et al. [20] found sparse pixel subsets suffice to attain high classification accuracy on popular  
121 image classification datasets, but evaluate interpretability methods rather than demonstrate spurious  
122 features or discover overinterpretation.
- 123 • Ghorbani et al. [21] introduce principles and methods for human-understandable concept-based  
124 explanations of ML models. In contrast, overinterpretation differs because the features we identify  
125 are semantically meaningless to humans, stem from single images, and are not aggregated into  
126 interpretable concepts. The existence of such subsets stemming from unmodified subsets of images  
127 suggests degeneracies in the underlying benchmark datasets and failures of modern CNN models  
128 to rely on more robust and interpretable signals in training datasets.
- 129 • Geirhos et al. [22] discuss the general problem of “shortcut learning” but do not recognize that  
130 5% (CIFAR-10) or 10% (ImageNet) spurious pixel-subsets are statistically valid signals in these  
131 datasets, nor characterize pixels that provide sufficient support and lead to overinterpretation.
- 132 • In natural language processing (NLP), Feng et al. [23] explored model pathologies using a similar  
133 technique, but did not analyze whether the semantically spurious patterns relied on are a statistical  
134 property of the dataset. Other work has demonstrated the presence of various spurious statistical  
135 shortcuts in major NLP benchmarks, showing this problem is not unique to CV [24].

## 136 3 Methods

### 137 3.1 Datasets and Models

138 CIFAR-10 [2] and ImageNet [3] have become two of the most popular image classification bench-  
139 marks. Most image classifiers are evaluated by the CV community based on their accuracy in one  
140 of these benchmarks. We also use the CIFAR-10-C dataset [25] to evaluate the extent to which our  
141 CIFAR-10 models can generalize to out-of-distribution (OOD) data. CIFAR-10-C contains variants  
142 of CIFAR-10 test images altered by various corruptions (e.g., Gaussian noise, motion blur). Where

143 computing sufficient input subsets on CIFAR-10-C images, we use a uniform random sample of 2000  
144 images across the entire CIFAR-10-C set. We use the ILSVRC2012 ImageNet dataset.

145 For CIFAR-10, we explore three common CNN architectures: a deep residual network with depth  
146 20 (ResNet20) [26], a v2 deep residual network with depth 18 (ResNet18) [27], and VGG16 [28].  
147 We train these networks using cross-entropy loss optimized via SGD with Nesterov momentum [29]  
148 and employ standard data augmentation strategies [27] (Section S1). After training many CIFAR-10  
149 networks individually, we construct four different ensemble classifiers by grouping various networks  
150 together. Each ensemble outputs the average prediction over its member networks (specifically,  
151 the arithmetic mean of their logits). For each of three architectures, we create a corresponding  
152 homogeneous ensemble by individually training five networks of that architecture. Each network  
153 has a different random initialization, which suffices to produce substantially different models despite  
154 having been trained on the same data [30]. Our fourth ensemble is heterogeneous, containing all 15  
155 networks (5 replicates of each of 3 distinct CNN architectures).

156 For ImageNet, we use a pre-trained Inception v3 model [31] that achieves 22.55% and 6.44% top-1  
157 and top-5 error [32].

### 158 3.2 Discovering Sufficient Features

159 **CIFAR-10.** We interpret the feature patterns learned by CIFAR-10 CNNs using the Sufficient  
160 Input Subsets (SIS) procedure [4], which produces rationales (SIS subsets) of a black-box model’s  
161 decision-making. SIS subsets are minimal subsets of input features (pixels) whose values alone  
162 suffice for the model to make the same decision as on the original input. Let  $f_c(x)$  denote the  
163 probability that an image  $x$  belongs to class  $c$ . An SIS subset  $S$  is a minimal subset of pixels of  $x$   
164 such that  $f_c(x_S) \geq \tau$ , where  $\tau$  is a prespecified confidence threshold and  $x_S$  is a modified input in  
165 which all information about values outside  $S$  are masked. We mask pixels by replacement with the  
166 mean value over all images (equal to zero when images have been normalized), which is presumably  
167 least informative to a trained classifier [4]. SIS subsets are found via a local backward selection  
168 algorithm applied to the function giving the confidence of the predicted (most likely) class.

169 **ImageNet.** We scale the SIS backward selection procedure to ImageNet with the introduction of  
170 Batched Gradient SIS, a gradient-based method to find sufficient input subsets on high-dimensional  
171 inputs. The sufficient input subsets discovered by Batched Gradient SIS are guaranteed to be sufficient,  
172 but may be larger than those discovered by the original exhaustive SIS algorithm. Here we find  
173 small SIS subsets with Batched Gradient SIS (Figure S10). Rather than separately masking every  
174 remaining pixel at each iteration to find the pixel whose masking least reduces  $f$ , we use the gradient  
175 of  $f$  with respect to the input pixels  $\mathbf{x}$  and mask  $M$ ,  $\nabla_M f(\mathbf{x} \odot (1 - M))$ , to order pixels (via a  
176 single backward pass). Instead of masking only one pixel per iteration, we mask larger subsets of  
177  $k \geq 1$  pixels per iteration. Given  $p$  input features, our Batched Gradient FindSIS procedure finds  
178 each SIS subset in  $\mathcal{O}(\frac{p}{k})$  evaluations of  $\nabla f$  (as opposed to  $\mathcal{O}(p^2)$  evaluations of  $f$  in FindSIS [4]).  
179 The complete Batched Gradient SIS algorithm is presented in Section S5.

### 180 3.3 Detecting Overinterpretation

181 We produce sparse variants of all train and test set images retaining 5% (CIFAR-10) or 10% (Im-  
182 ageNet) of pixels in each image. Our goal is to identify sparse pixel-subsets that contain feature  
183 patterns the model identifies as strong class-evidence as it classifies an image. We identify pixels  
184 to retain based on sorting by SIS BackSelect [4] (CIFAR-10) or our Batched Gradient BackSelect  
185 procedure (ImageNet). These backward selection (BS) pixel-subset images contain the final pixels  
186 (with their same RGB values as in the original images) while all other pixels’ values are replaced with  
187 zero. Note that we apply backward selection to the function giving the confidence of the *predicted*  
188 class from the original model to prevent adding information about the true class for misclassified  
189 images, and we use the true labels for training/evaluating models on pixel-subsets. As backward  
190 selection is applied locally on each image, the specific pixels retained differ across images.

191 We train new classifiers on solely these pixel-subsets of training images and evaluate accuracy on  
192 corresponding pixel-subsets of test images to determine whether such pixel-subsets are statistically  
193 valid for generalization in the benchmark. We use the same training setup and hyperparameters  
194 (Section 3.1) without data augmentation of training images (results with data augmentation in

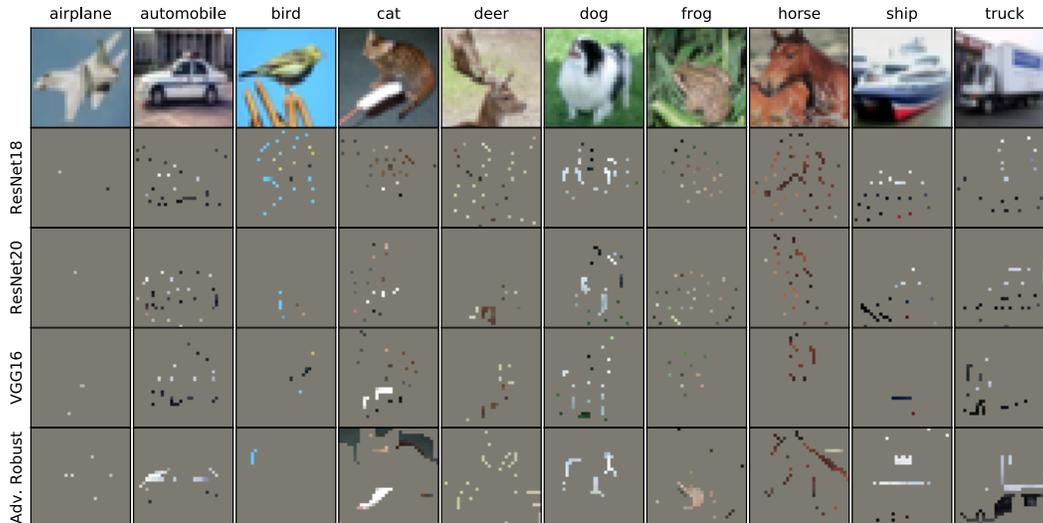


Figure 1: Sufficient input subsets (SIS) for a sample of CIFAR-10 test images (top). Each SIS image shown below is classified by the respective model with  $\geq 99\%$  confidence.

195 Table S1). We consider a model to overinterpret its input when these signals can generalize to test  
 196 data but lack semantic meaning (Section 3.4).

### 197 3.4 Human Classification Benchmark

198 To evaluate whether sparse pixel-subsets of images can be accurately classified by humans, we asked  
 199 four participants to classify images containing various degrees of masking. We randomly sampled  
 200 100 images from the CIFAR-10 test set (10 images per class) that were correctly and confidently  
 201 ( $\geq 99\%$  confidence) classified by our models, and for each image, kept only 5%, 30%, or 50% of  
 202 pixels as ranked by backward selection (all other pixels masked). Backward selection image subsets  
 203 are sampled across our three models. Since larger subsets of pixels are by construction supersets  
 204 of smaller subsets identified by the same model, we presented each batch of 100 images in order  
 205 of increasing subset size and shuffled the order of images within each batch. Users were asked to  
 206 classify each of the 300 images as one of the 10 classes in CIFAR-10 and were not provided training  
 207 images. The same task was given to each user (and is shown in Section S4).

## 208 4 Results

### 209 4.1 CNNs Classify Images Using Spurious Features

210 **CIFAR-10.** Figure 1 shows example SIS subsets (threshold 0.99) from CIFAR-10 test images  
 211 (additional examples in Section S2). These SIS subset images are confidently and correctly classified  
 212 by each model with  $\geq 99\%$  confidence toward the predicted class. We observe these SIS subsets  
 213 are highly sparse and the average SIS size at this threshold is  $< 5\%$  of each image (see Figure 5),  
 214 suggesting these CNNs confidently classify images that appear nonsensical to humans (Section 4.3),  
 215 leading to concern about their robustness and generalizability.

216 We retain 5% of pixels in each image using local backward selection and mask the remaining  
 217 95% with zeros (Section 3.3) and find models trained on full images classify these pixel-subsets as  
 218 accurately as full images (Table 1). Figure 2a shows the pixel locations and confidence of these 5%  
 219 pixel-subsets across all CIFAR-10 test images. Moreover, the CNNs are more confident on these  
 220 pixels subsets than on full images: the mean drop in confidence for the predicted class between  
 221 original images and these 5% subsets is  $-0.035$  (std dev. = 0.107),  $-0.016$  (0.094), and  $-0.012$   
 222 (0.074) computed over all CIFAR-10 test images for our ResNet20, ResNet18, and VGG16 models,  
 223 respectively, suggesting severe overinterpretation (negative values imply greater confidence on the

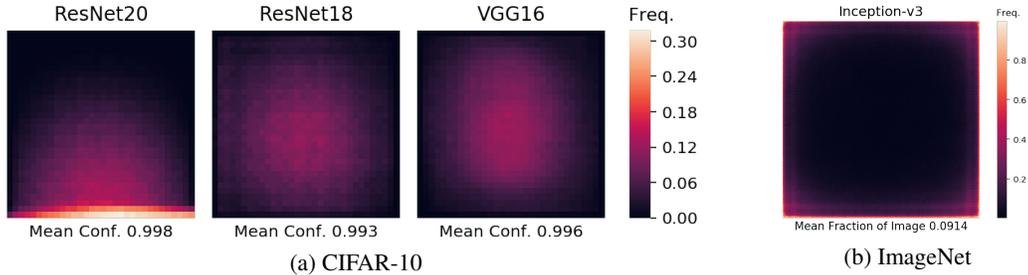


Figure 2: Heatmaps of pixel locations comprising pixel-subsets. Frequency indicates fraction of subsets containing each pixel. (a) 5% pixel-subsets across CIFAR-10 test set for each model. Mean confidence indicates confidence on 5% pixel-subsets. (b) Sufficient input subsets (threshold 0.9) across ImageNet validation images from Inception v3.



Figure 3: Sufficient input subsets (threshold 0.9) for example ImageNet validation images. The bottom row shows the corresponding images with all pixels outside of each SIS subset masked but are still classified by the Inception v3 model with  $\geq 90\%$  confidence.

224 5% subsets). We find pixel-subsets chosen via backward selection are significantly more predictive  
 225 than equally large pixel-subsets chosen uniformly at random from each image (Table 1).

226 We also find SIS subsets confidently classified by one model do not transfer to other models. For  
 227 instance, 5% pixel-subsets derived from CIFAR-10 test images using one ResNet18 model (which  
 228 classifies them with 94.8% accuracy) are only classified with 25.8%, 29.2%, and 27.5% accuracy by  
 229 another ResNet18 replicate, ResNet20, and VGG16 models, respectively, suggesting there exist many  
 230 different statistical patterns that a flexible model might learn to rely on, and thus CIFAR-10 image  
 231 classification remains a highly underdetermined problem. Training classifiers that make predictions  
 232 for the right reasons may require clever regularization strategies and architecture design to ensure  
 233 models favor salient features over spurious pixel subsets.

234 While recent work has suggested semantics can be better captured by models that are robust to  
 235 adversarial inputs that fool standard neural networks via human-imperceptible modifications to  
 236 images [19, 33], we explore a wide residual network that is adversarially robust for CIFAR-10  
 237 classification [19] and find evidence of overinterpretation (Figure 1). This finding suggests adversarial  
 238 robustness alone does not prevent models from overinterpreting spurious signals in CIFAR-10.

239 **ImageNet.** We also find models trained on ImageNet images suffer from severe overinterpretation.  
 240 Figure 3 shows example SIS subsets (threshold 0.9) found via Batched Gradient SIS on images  
 241 confidently classified by the pre-trained Inception v3 (additional examples in Figures S8 and S9).  
 242 These SIS subsets appear visually nonsensical, yet the network classifies them with  $\geq 90\%$  confidence.  
 243 We find SIS pixels are concentrated outside of the actual object that determines the class label. For  
 244 example, in the “pizza” image, the SIS is concentrated on the shape of the plate and the background  
 245 table, rather than the pizza itself, suggesting the model could generalize poorly on images containing  
 246 different circular items on a table. In the “giant panda” image, the SIS contains bamboo, which  
 247 likely appeared in the collection of ImageNet photos for this class. In the “traffic light” and “street

Table 1: Accuracy of CIFAR-10 classifiers trained and evaluated on full images, 5% backward selection (BS) pixel-subsets, and 5% random pixel-subsets. Where possible, we report accuracy as mean  $\pm$  standard deviation (%) over five runs. For training on BS subsets, we run BS on all images for a single model of each type and average over five models trained on these subsets.

Model	Train On	Evaluate On	CIFAR-10 Test Acc.	CIFAR-10-C Acc.
ResNet20	Full Images	Full Images	92.52 $\pm$ 0.09	69.44 $\pm$ 0.52
		5% BS Subsets	92.48	70.65
		5% Random	9.98 $\pm$ 0.03	10.02 $\pm$ 0.01
	5% BS Subsets	5% BS Subsets	92.49 $\pm$ 0.02	70.58 $\pm$ 0.03
	5% Random	5% Random	50.25 $\pm$ 0.19	44.04 $\pm$ 0.33
	Input Dropout (Full)	Input Dropout (Full)	91.02 $\pm$ 0.25	75.46 $\pm$ 0.74
ResNet18	Full Images	Full Images	95.17 $\pm$ 0.21	75.08 $\pm$ 0.20
		5% BS Subsets	94.76	75.15
		5% Random	10.08 $\pm$ 0.15	10.08 $\pm$ 0.07
	5% BS Subsets	5% BS Subsets	94.96 $\pm$ 0.04	75.25 $\pm$ 0.05
	5% Random	5% Random	51.27 $\pm$ 0.82	45.24 $\pm$ 0.45
	Input Dropout (Full)	Input Dropout (Full)	94.15 $\pm$ 0.26	80.35 $\pm$ 0.39
VGG16	Full Images	Full Images	93.69 $\pm$ 0.12	74.14 $\pm$ 0.45
		5% BS Subsets	93.27	73.95
		5% Random	10.02 $\pm$ 0.18	9.97 $\pm$ 0.18
	5% BS Subsets	5% BS Subsets	92.60 $\pm$ 0.08	73.27 $\pm$ 0.18
	5% Random	5% Random	53.66 $\pm$ 1.96	46.88 $\pm$ 1.27
	Input Dropout (Full)	Input Dropout (Full)	91.09 $\pm$ 0.15	80.43 $\pm$ 0.24
Ensemble (ResNet18)	Full Images	Full Images	96.07	77.00
		5% Random	9.98	10.01

248 sign” images, the SIS consists of pixels in sky, suggesting that autonomous vehicle systems that may  
 249 depend on these models should be carefully evaluated for overinterpretation pathologies.

250 Figure 2b shows SIS pixel locations from a random sample of 1000 ImageNet validation images. We  
 251 find concentration along image borders, suggesting the model relies heavily on image backgrounds  
 252 and suffers from severe overinterpretation. This is a serious problem as objects determining ImageNet  
 253 classes are often located near image centers, and thus this network fails to focus on salient features.

## 254 4.2 Sparse Subsets are Real Statistical Patterns

255 The overconfidence of CNNs for image classification [34] may lead one to wonder whether the  
 256 observed overconfidence on semantically meaningless SIS subsets is an artifact of calibration rather  
 257 than true statistical signals in the dataset. We train models on 5% pixel-subsets of CIFAR-10 training  
 258 images found via backward selection (Section 3.3). We find models trained solely on these pixel-  
 259 subsets can classify corresponding test image pixel-subsets with minimal accuracy loss compared to  
 260 models trained on full images (Table 1). As a baseline to the 5% pixel-subsets identified by backward  
 261 selection, we create variants of all images where the 5% pixel-subsets are selected at random from  
 262 each image (rather than by backward selection) and use the same random pixel-subsets for training  
 263 each new model. Models trained on random subsets have significantly lower test accuracy compared  
 264 to models trained on 5% pixel-subsets from backward selection (Table 1). We observe, however, that  
 265 random 5% subsets of images still capture enough signal to predict roughly 5 times better than blind  
 266 guessing, but do not capture nearly enough information for models to make accurate predictions.

267 We found that the 5% backward selection pixel-subsets did not contain model specific features,  
 268 and thus reflected valid predictive signals regardless of the model architecture employed for subset  
 269 discovery. Our hypothesis was that 5% pixel-subsets discovered with one architecture would provide  
 270 robust performance when used to train and evaluate a second architecture. We found this hypothesis

271 supported for all six pairs of subset discovery and train-test architectures evaluated (Table S2). These  
272 results demonstrate that the highly sparse subsets found via backward selection offer a valid predictive  
273 signal in the CIFAR-10 benchmark exploited by models to attain high test accuracy.

274 We observe similar results on ImageNet. Inception v3 trained on 10% pixel-subsets of ImageNet  
275 training images achieves 71.4% accuracy (mean over 5 runs) on the corresponding pixel-subset  
276 ImageNet validation set (Table S4). Additional ImageNet results for Inception v3 and ResNet-50,  
277 including training and evaluation on random subsets, are provided in Table S4.

### 278 4.3 Humans Struggle to Classify Sparse Subsets

279 We find a strong correlation between the fraction of unmasked pixels in each image and human  
280 classification accuracy ( $R^2 = 0.94$ , Figure S7). Human accuracy on 5% pixel-subsets of CIFAR-10  
281 images (mean = 19.2%, std dev = 4.8%, Table S3) is significantly lower than on original, unmasked  
282 images (roughly 94% [35]), though greater than random guessing, presumably due to correlations  
283 between labels and features such as color (e.g., blue sky suggests airplane, ship, or bird).

284 However, CNNs (even when trained on full images and achieve accuracy on par with human accuracy  
285 on full images) classify these sparse image subsets with very high accuracy (Table I), indicating  
286 benchmark images contain statistical signals that are not salient to humans. Models solely trained  
287 to minimize prediction error may thus latch onto these signals while still accurately generalizing to  
288 test data, but may behave counterintuitively when fed images from a different source that does not  
289 share these exact statistics. The strong correlation between the size of CIFAR-10 pixel-subsets and  
290 the corresponding human classification accuracy suggests larger subsets contain more semantically  
291 salient content. Thus, a model whose decisions have larger corresponding SIS subsets presumably  
292 exhibits less overinterpretation than one with smaller SIS subsets, as we investigate in Section 4.4.

### 293 4.4 SIS Size is Related to Model Accuracy

294 Given that smaller SIS contain fewer salient features according to human classifiers, models that  
295 justify their classifications based on sparse SIS subsets may be limited in terms of attainable accuracy,  
296 particularly in out-of-distribution settings. Here, we investigate the relationship between a single  
297 model’s predictive accuracy and the size of the SIS subsets in which it identifies class-evidence. We  
298 draw no conclusions between models as they are uncalibrated. For each of our three classifiers, we  
299 compute the average SIS size increase for correctly classified images as compared to incorrectly  
300 classified images (expressed as a percentage). We find SIS subsets of correctly classified images are  
301 consistently significantly larger than those of misclassified images at all SIS confidence thresholds for  
302 both CIFAR-10 test images (Figure 4) and CIFAR-10-C OOD images (Figure S3). This is especially  
303 striking given model confidence is uniformly lower on the misclassified inputs (Figure S4). Lower  
304 confidence would normally imply a larger SIS subset at a given confidence level, as one expects  
305 fewer pixels can be masked before the model’s confidence drops below the SIS threshold. Thus, we  
306 can rule out overall model confidence as an explanation of the smaller SIS of misclassified images.  
307 This result suggests the sparse SIS subsets highlighted in this paper are not just a curiosity, but may  
308 be leading to poor generalization on real images.

### 309 4.5 Mitigating Overinterpretation

310 **Ensembling.** Model ensembling is known to improve classification performance [36, 37]. As we  
311 found pixel-subset size to be strongly correlated with human pixel-subset classification accuracy  
312 (Section 4.3), our metric for measuring how much ensembling may alleviate overinterpretation is the  
313 increase in SIS subset size. We find ensembling uniformly increases test accuracy as expected but  
314 also increases the SIS size (Figure 5), hence mitigating overinterpretation.

315 We conjecture the cause of both the increase in the accuracy and SIS size for ensembles is the same.  
316 We observe that SIS subsets are generally not transferable from one model to another — i.e., an SIS  
317 for one model is rarely an SIS for another (Section 4.1). Thus, different models rely on different  
318 independent signals to arrive at the same prediction. An ensemble bases its prediction on multiple  
319 such signals, increasing predictive accuracy and SIS subset size by requiring simultaneous activation  
320 of multiple independently trained feature detectors. We find SIS subsets of the ensemble are larger  
321 than the SIS of its individual members (examples in Figure S2).

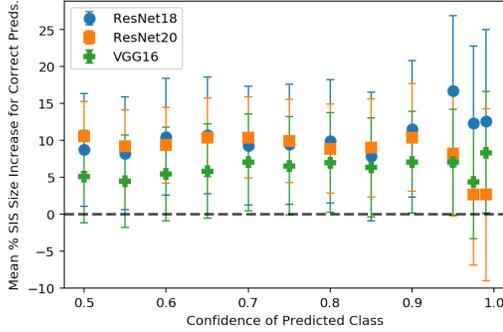


Figure 4: Percentage increase in mean SIS size of correctly classified compared to misclassified CIFAR-10 test images. Positive values indicate larger mean SIS size for correctly classified images. Error bars indicate 95% confidence interval for the difference in means.

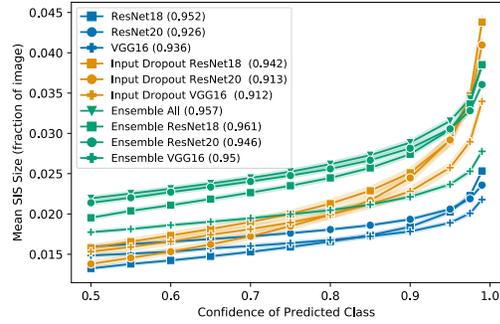


Figure 5: Mean SIS size on CIFAR-10 test images as SIS threshold varies. SIS size indicates fraction of pixels necessary for model to make the same prediction at each confidence threshold. Model accuracies are shown in the legend. 95% confidence intervals are shaded around each mean.

322 **Input Dropout.** We apply input dropout [38] to both train and test images. We retain each input  
 323 pixel with probability  $p = 0.8$  and set the values of dropped pixels to zero. We find a small decrease  
 324 in CIFAR-10 test accuracy for models regularized with input dropout though find a significant ( $\sim 6\%$ )  
 325 increase in OOD test accuracy on CIFAR-10-C images (Table 1, Figure S5). Figure 5 shows a  
 326 corresponding increase in SIS subset size for these models, suggesting input dropout applied at train  
 327 and test time helps to mitigate overinterpretation. We conjecture that random dropout of input pixels  
 328 disrupts spurious signals that lead to overinterpretation.

## 329 5 Discussion

330 We find that modern image classifiers overinterpret small nonsensical patterns present in popular  
 331 benchmark datasets, identifying strong class evidence in the pixel-subsets that constitute these patterns.  
 332 We introduced the Batched Gradient SIS method for the efficient discovery of such patterns. Despite  
 333 their lack of salient features, these sparse pixel-subsets are underlying statistical signals that suffice  
 334 to accurately generalize from the benchmark training data to the benchmark test data. We found that  
 335 different models rationalize their predictions based on different sufficient input subsets, suggesting  
 336 optimal image classification rules remain highly underdetermined by the training data. In high-stakes  
 337 applications, we recommend ensembles of networks or regularization via input dropout.

338 Our results call into question model interpretability methods whose outputs are encouraged to align  
 339 with prior human beliefs of proper classifier operating behavior [39]. Given the existence of non-  
 340 salient pixel-subsets that alone suffice for correct classification, a model may solely rely on such  
 341 patterns. In this case, an interpretability method that faithfully describes the model should output  
 342 these nonsensical rationales, whereas interpretability methods that bias rationales toward human  
 343 priors may produce results that mislead users to think their models behave as intended.

344 Mitigating overinterpretation and the broader task of ensuring classifiers are accurate for the right  
 345 reasons remain significant challenges for ML. While we identify strategies for partially mitigating  
 346 overinterpretation, additional research needs to develop ML methods that rely exclusively on well-  
 347 formed interpretable inputs, and methods for creating training data that do not contain spurious  
 348 signals. One alternative is to regularize CNNs by constraining the pixel attributions generated via  
 349 a saliency map [40, 42]. Unfortunately, such methods require a human annotator to highlight the  
 350 correct pixels as an auxiliary supervision signal. Saliency maps have also been shown to provide  
 351 unreliable insights into model operating behavior and must be interpreted as approximations [43].  
 352 In contrast, our SIS subsets constitute actual pathological examples that have been misconstrued by  
 353 the model. An important application of our methods is the evaluation of training datasets to ensure  
 354 decisions are made on interpretable rather than spurious signals. We found popular image datasets  
 355 contain such spurious signals, and the resulting overinterpretation may be difficult to overcome with  
 356 ML methods alone. We provide an open-source implementation of our methods.

## 357 References

- 358 [1] Neel V. Patel. *Why Doctors Aren't Afraid of Better, More Efficient*  
359 *AI Diagnosing Cancer*, 2017. URL [https://www.thedailybeast.com/](https://www.thedailybeast.com/why-doctors-arent-afraid-of-better-more-efficient-ai-diagnosing-cancer)  
360 [why-doctors-arent-afraid-of-better-more-efficient-ai-diagnosing-cancer](https://www.thedailybeast.com/why-doctors-arent-afraid-of-better-more-efficient-ai-diagnosing-cancer). Accessed  
361 September 27, 2020.
- 362 [2] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of  
363 Toronto, 2009.
- 364 [3] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,  
365 Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large  
366 Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252,  
367 2015. doi: 10.1007/s11263-015-0816-y.
- 368 [4] Brandon Carter, Jonas Mueller, Siddhartha Jain, and David Gifford. What made you do this? Understanding  
369 black-box decisions with sufficient input subsets. In *The 22nd International Conference on Artificial*  
370 *Intelligence and Statistics*, pages 567–576, 2019.
- 371 [5] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry.  
372 Adversarial examples are not bugs, they are features. In *Advances in Neural Information Processing*  
373 *Systems*, pages 125–136, 2019.
- 374 [6] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528.  
375 IEEE, 2011.
- 376 [7] Tatiana Tommasi, Novi Patricia, Barbara Caputo, and Tinne Tuytelaars. A deeper look at dataset bias. In  
377 *Domain adaptation in computer vision applications*, pages 37–55. Springer, 2017.
- 378 [8] Amir Rosenfeld, Richard Zemel, and John K Tsotsos. The elephant in the room. *arXiv preprint*  
379 *arXiv:1808.03305*, 2018.
- 380 [9] Rakshith Shetty, Bernt Schiele, and Mario Fritz. Not using the car to see the sidewalk—quantifying and  
381 controlling the effects of context in classification and segmentation. In *Proceedings of the IEEE/CVF*  
382 *Conference on Computer Vision and Pattern Recognition*, pages 8218–8226, 2019.
- 383 [10] Krishna Kumar Singh, Dhruv Mahajan, Kristen Grauman, Yong Jae Lee, Matt Feiszli, and Deepti Ghadi-  
384 yaram. Don't judge an object by its context: Learning to overcome contextual bias. In *Proceedings of the*  
385 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11070–11078, 2020.
- 386 [11] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and  
387 Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really learn. *Nature*  
388 *Communications*, 10(1):1–8, 2019.
- 389 [12] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Texture and art with deep neural networks.  
390 *Current Opinion in Neurobiology*, 46:178–186, 2017.
- 391 [13] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland  
392 Brendel. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and  
393 robustness. *arXiv preprint arXiv:1811.12231*, 2018.
- 394 [14] Wieland Brendel and Matthias Bethge. Approximating CNNs with Bag-of-Local-Features models works  
395 surprisingly well on ImageNet. *Proceedings of the International Conference on Learning Representations*,  
396 2019.
- 397 [15] Robert Geirhos, Carlos R Medina Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A  
398 Wichmann. Generalisation in humans and deep neural networks. *arXiv preprint arXiv:1808.08750*, 2018.
- 399 [16] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence  
400 predictions for unrecognizable images. In *CVPR*, 2015.
- 401 [17] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples.  
402 *arXiv preprint arXiv:1412.6572*, 2014.
- 403 [18] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence  
404 predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and*  
405 *pattern recognition*, pages 427–436, 2015.
- 406 [19] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards  
407 deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

- 408 [20] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability  
409 methods in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 9734–  
410 9745, 2019.
- 411 [21] Amirata Ghorbani, James Wexler, James Zou, and Been Kim. Towards automatic concept-based explana-  
412 tions. *arXiv preprint arXiv:1902.03129*, 2019.
- 413 [22] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias  
414 Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*,  
415 2(11):665–673, 2020.
- 416 [23] Shi Feng, Eric Wallace, II Grissom, Mohit Iyyer, Pedro Rodriguez, Jordan Boyd-Graber, et al. Pathologies  
417 of neural models make interpretations difficult. *arXiv preprint arXiv:1804.07781*, 2018.
- 418 [24] Timothy Niven and Hung-Yu Kao. Probing neural network comprehension of natural language arguments.  
419 *ACL*, 2019.
- 420 [25] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions  
421 and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.
- 422 [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
423 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 424 [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks.  
425 In *European conference on computer vision*, pages 630–645. Springer, 2016.
- 426 [28] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recogni-  
427 tion. *arXiv preprint arXiv:1409.1556*, 2014.
- 428 [29] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and  
429 momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- 430 [30] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped  
431 dqn. In *Advances in neural information processing systems*, pages 4026–4034, 2016.
- 432 [31] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the  
433 inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and  
434 pattern recognition*, pages 2818–2826, 2016.
- 435 [32] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,  
436 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style, high-performance deep  
437 learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019.
- 438 [33] Shibani Santurkar, Andrew Ilyas, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry.  
439 Image synthesis with a single (robust) classifier. In *Advances in Neural Information Processing Systems*,  
440 pages 1260–1271, 2019.
- 441 [34] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks.  
442 In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1321–1330.  
443 JMLR. org, 2017.
- 444 [35] Andrej Karpathy. Lessons learned from manually classifying CIFAR-10. *Published online at  
445 <http://karpathy.github.io/2011/04/27/manually-classifying-cifar10>*, 2011.
- 446 [36] King-Shy Goh, Edward Chang, and Kwang-Ting Cheng. SVM binary classifier ensembles for image  
447 classification. In *Proceedings of the tenth international conference on Information and knowledge  
448 management*, pages 395–402. ACM, 2001.
- 449 [37] Cheng Ju, Aurélien Bibaut, and Mark van der Laan. The relative performance of ensemble methods  
450 with deep convolutional neural networks for image classification. *Journal of Applied Statistics*, 45(15):  
451 2800–2818, 2018.
- 452 [38] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout:  
453 A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:  
454 1929–1958, 2014.
- 455 [39] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity  
456 checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515, 2018.

- 457 [40] Andrew Slavin Ross, Michael C Hughes, and Finale Doshi-Velez. Right for the right reasons: training  
 458 differentiable models by constraining their explanations. In *Proceedings of the 26th International Joint*  
 459 *Conference on Artificial Intelligence*, pages 2662–2670, 2017.
- 460 [41] Becks Simpson, Francis Dutil, Yoshua Bengio, and Joseph Paul Cohen. GradMask: Reduce overfitting by  
 461 regularizing saliency. *arXiv preprint arXiv:1904.07478*, 2019.
- 462 [42] Joseph D Viviano, Becks Simpson, Francis Dutil, Yoshua Bengio, and Joseph Paul Cohen. Underwhelming  
 463 generalization improvements from controlling feature attribution. *arXiv preprint arXiv:1910.00199*, 2019.
- 464 [43] Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne,  
 465 Dumitru Erhan, and Been Kim. The (un) reliability of saliency methods. In *Explainable AI: Interpreting,*  
 466 *Explaining and Visualizing Deep Learning*, pages 267–280. Springer, 2019.

## 467 Checklist

- 468 1. For all authors...
- 469 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 470 contributions and scope? [Yes]
- 471 (b) Did you describe the limitations of your work? [Yes] We demonstrate that ensembling  
 472 and input dropout (Section 4.5) mitigate but do not completely prevent overinterpretation as overinterpretation is caused by spurious statistical signals in training data (discussed in Section 5).
- 473 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We  
 474 discuss implications for dataset curation in Section 5. One potential consequence of  
 475 this work is that training datasets may become more complex and costly to generate to  
 476 remove the kinds of degenerate signals we have observed.
- 477 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 478 them? [Yes]
- 479
- 480 2. If you are including theoretical results...
- 481 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 482 (b) Did you include complete proofs of all theoretical results? [N/A]
- 483
- 484 3. If you ran experiments...
- 485 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
 486 perimental results (either in the supplemental material or as a URL)? [Yes] See  
 487 supplementary material.
- 488 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 489 were chosen)? [Yes] See Sections 3.1 and S1 (model training), Section 3.2 (SIS), and  
 490 Sections 3.3 and S3 (overinterpretation).
- 491 (c) Did you report error bars (e.g., with respect to the random seed after running exper-  
 492 iments multiple times)? [Yes] Models were trained multiple times with different  
 493 random seeds, and accuracies in Table 1 are reported as mean  $\pm$  standard deviation.  
 494 Figures 4 and 5 show error bars indicating 95% confidence intervals.
- 495 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 496 of GPUs, internal cluster, or cloud provider)? [Yes] See Section S1.
- 497
- 498 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 499 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section S1
- 500 (b) Did you mention the license of the assets? [N/A] We used the CIFAR-10 and ImageNet  
 501 datasets, and could not locate specific license information.
- 502 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
 503 Our code is included in the supplemental material and will be released on GitHub  
 504 under an open-source license.
- 505 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
 506 using/curating? [N/A] Previously published data were utilized.
- 507 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
 508 information or offensive content? [N/A] Previously published data were utilized.

- 508 5. If you used crowdsourcing or conducted research with human subjects...
- 509 (a) Did you include the full text of instructions given to participants and screenshots, if
- 510 applicable? [Yes] See Sections 3.4 and S4 and Figure S6.
- 511 (b) Did you describe any potential participant risks, with links to Institutional Review
- 512 Board (IRB) approvals, if applicable? [N/A] IRB approval was not required.
- 513 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 514 spent on participant compensation? [N/A] Users were volunteers.