

# Overinterpretation reveals image classification model pathologies

**Brandon Carter**  
MIT CSAIL  
bcarter@csail.mit.edu

**Siddhartha Jain**  
MIT CSAIL

**Jonas Mueller**  
Amazon Web Services

**David Gifford**  
MIT CSAIL  
gifford@mit.edu

## Abstract

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

*Adversarial examples are generated by modifying or shuffling the image pixels. While in overinterpretation relies on masking the image pixels without any shuffling.*

## 1 Introduction

Well-founded decisions by machine learning (ML) systems are critical for high-stakes applications such as autonomous vehicles and medical diagnosis. Pathologies in models and their respective training datasets can result in unintended behavior during deployment if the systems are confronted with novel situations. For example, a medical image classifier for cancer detection attained high accuracy in benchmark test data, but was found to base decisions upon presence of rulers in an image (present when dermatologists already suspected cancer) [1]. We define model *overinterpretation* to occur when a classifier finds strong class-evidence in regions of an image that contain no semantically salient features. Overinterpretation is related to overfitting, but overfitting can be diagnosed via reduced test accuracy. Overinterpretation can stem from true statistical signals in the underlying dataset distribution that happen to arise from particular properties of the data source (e.g., dermatologists' rulers). Thus, overinterpretation can be harder to diagnose as it admits decisions that are made by statistically valid criteria, and models that use such criteria can excel at benchmarks. We demonstrate overinterpretation occurs with unmodified subsets of the original images. In contrast to *adversarial examples* that modify images with extra information, overinterpretation is based on real patterns already present in the training data that also generalize to the test distribution. Hidden statistical signals of benchmark datasets can result in models that overinterpret or do not generalize to new data from a different distribution. Computer vision (CV) research relies on datasets like CIFAR-10 [2] and ImageNet [3] to provide standardized performance benchmarks. Here, we analyze the overinterpretation of popular CNN architectures on these benchmarks to characterize pathologies.

*Overinterpretation arises from the true patterns present in the dataset itself.*

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

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

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

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

Code for this paper is available at: <https://github.com/gifford-lab/overinterpretation>.

## 2 Related Work

While existing work has demonstrated numerous distinct flaws in deep image classifiers our paper demonstrates a new distinct flaw, overinterpretation, previously undocumented in the literature. There has been substantial research on understanding dataset bias in CV [6, 7] and the fragility of image classifiers deployed outside benchmark settings. We extend previous work on sufficient input subsets (SIS) [4] with the Batched Gradient SIS method, and use this method to show that ImageNet sufficient input subset pixels for training and testing often exist at image borders. Many alternative interpretability methods also aim to understand models by extracting *rationales* (pixel-subsets) that

Model explainability can be used to account for the features which accounts most for the final output. The author used the concept of SIS which is basically the minimum number of pixel required in image which yield same performance as when used with all pixels.

Overinterpretation shows that classifier is learning inherent shortcut pattern instead of learning meaningful complicated patterns

Adversarial models also suffers form overinterpretation.

The author defines how their work are different from the previous work which also presented flaws in classifier

The major difference are this

1. In this paper they used SIS which take pixels which are not interpretable by humans.Ã€

2. Imgaes have not been modified it is masked.

3. They exposed though classifiers learn spurious but statistically valid signal.

provide positive evidence for a class [8–11], and we adopt SIS throughout this work as a particularly straightforward method for producing such rationales. This prior work (including SIS [4]) is limited to understanding models and does not use the enhanced understanding of models to identify the overinterpretation flaw discovered in this paper. We contrast the issue of overinterpretation against other previously known model flaws below:

- Image classifiers have been shown to be fragile when objects from one image are transplanted in another image [12], and can be biased by object context [13, 14]. In contrast, overinterpretation differs because we demonstrate that highly sparse, unmodified subsets of pixels in images suffice for image classifiers to make the same predictions as on the full images.
- Lapuschkin et al. [15] demonstrate that DNNs can learn to rely on spurious signals in datasets, including source tags and artificial padding, but which are still human-interpretable. In contrast, the patterns we identify are minimal collections of pixels in images that are semantically meaningless to humans (they do not comprise human-interpretable parts of images). We demonstrate such patterns generalize to the test distribution suggesting they arise from degenerate signals in popular benchmarks, and thus models trained on these datasets may fail to generalize to real-world data.
- CNNs in particular have been conjectured to pick up on localized features like texture instead of more global features like object shape [16, 17]. Brendel and Bethge [18] show CNNs trained on natural ImageNet images may rely on local features and, unlike humans, are able to classify texturized images, suggesting ImageNet alone is insufficient to force DNNs to rely on more causal representations. Our work demonstrates another source of degeneracy of popular image datasets, where sparse, unmodified subsets of training images that are meaningless to humans can enable a model to generalize to test data. We provide one explanation for why ImageNet-trained models may struggle to generalize to out-of-distribution data.
- Geirhos et al. [19] discover that DNNs trained on distorted images fail to generalize as well as human observers when trained under image distortions. In contrast, overinterpretation reveals a different failure mode of DNNs, whereby models latch onto spurious but statistically valid sets of features in undistorted images. This phenomenon can limit the ability of a DNN to generalize to real-world data even when trained on natural images.
- Other work has shown deep image classifiers can make confident predictions on nonsensical patterns [20], and the susceptibility of DNNs to adversarial examples or synthetic images has been widely studied [5, 21–23]. However, these adversarial examples synthesize artificial images or modify real images with auxiliary information. In contrast, we demonstrate overinterpretation of unmodified subsets of actual training images, indicating the patterns are already present in the original dataset. We further demonstrate that such signals in training data actually generalize to the test distribution and that adversarially robust models also suffer from overinterpretation.
- Hooker et al. [24] found sparse pixel subsets suffice to attain high classification accuracy on popular image classification datasets, but evaluate interpretability methods rather than demonstrate spurious features or discover overinterpretation.
- Ghorbani et al. [25] introduce principles and methods for human-understandable concept-based explanations of ML models. In contrast, overinterpretation differs because the features we identify are semantically meaningless to humans, stem from single images, and are not aggregated into interpretable concepts. The existence of such subsets stemming from unmodified subsets of images suggests degeneracies in the underlying benchmark datasets and failures of modern CNN models to rely on more robust and interpretable signals in training datasets.
- Geirhos et al. [26] discuss the general problem of “shortcut learning” but do not recognize that 5% (CIFAR-10) or 10% (ImageNet) spurious pixel-subsets are statistically valid signals in these datasets, nor characterize pixels that provide sufficient support and lead to overinterpretation.
- In natural language processing (NLP), Feng et al. [27] explored model pathologies using a similar technique, but did not analyze whether the semantically spurious patterns relied on are a statistical property of the dataset. Other work has demonstrated the presence of various spurious statistical shortcuts in major NLP benchmarks, showing this problem is not unique to CV [28].

### 3 Methods

#### 3.1 Datasets and Models

CIFAR-10 [2] and ImageNet [3] have become two of the most popular image classification benchmarks. Most image classifiers are evaluated by the CV community based on their accuracy in one of these benchmarks. We also use the CIFAR-10-C dataset [29] to evaluate the extent to which our CIFAR-10 models can generalize to out-of-distribution (OOD) data. CIFAR-10-C contains variants of CIFAR-10 test images altered by various corruptions (e.g., Gaussian noise, motion blur). When computing sufficient input subsets on CIFAR-10-C images, we use a uniform random sample of 2000 images across the entire CIFAR-10-C set. Additional results on CIFAR-10.1 v6 [30] are presented in Table S4. We use the ILSVRC2012 ImageNet dataset [3].

For CIFAR-10, we explore three common CNN architectures: a deep residual network with depth 20 (ResNet20) [31], a v2 deep residual network with depth 18 (ResNet18) [32], and VGG16 [33]. We train these networks using cross-entropy loss optimized via SGD with Nesterov momentum [34] and employ standard data augmentation strategies [32] (Section S2). After training many CIFAR-10 networks individually, we construct four different ensemble classifiers by grouping various networks together. Each ensemble outputs the average prediction over its member networks (specifically, the arithmetic mean of their logits). For each of three architectures, we create a corresponding homogeneous ensemble by individually training five networks of that architecture. Each network has a different random initialization, which suffices to produce substantially different models despite having been trained on the same data [35]. Our fourth ensemble is heterogeneous, containing all 15 networks (5 replicates of each of 3 distinct CNN architectures).

For ImageNet, we use a pre-trained Inception v3 model [36] that achieves 22.55% and 6.44% top-1 and top-5 error [37]. Additional results from an ImageNet ResNet50 are presented in Section S6.

#### 3.2 Discovering Sufficient Features

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

**ImageNet.** We scale the SIS backward selection procedure to ImageNet with the introduction of Batched Gradient SIS, a gradient-based method to find sufficient input subsets on high-dimensional inputs. The sufficient input subsets discovered by Batched Gradient SIS are guaranteed to be sufficient, but may be larger than those discovered by the original exhaustive SIS algorithm. Here we find small SIS subsets with Batched Gradient SIS (Figure S15). Rather than separately masking every remaining pixel at each iteration to find the pixel whose masking least reduces  $f$ , we use the gradient of  $f$  with respect to the input pixels  $x$  and mask  $M$ ,  $\nabla_M f(x \odot (1 - M))$ , to order pixels (via a single backward pass). Instead of masking only one pixel per iteration, we mask larger subsets of  $k \geq 1$  pixels per iteration. Given  $p$  input features, our Batched Gradient FindSIS procedure finds 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]). The complete Batched Gradient SIS algorithm is presented in Section S1.

Author used  
batch gradient  
SIS to scale it to  
ImageNet.

#### 3.3 Detecting Overinterpretation

We produce sparse variants of all train and test set images retaining 5% (CIFAR-10) or 10% (ImageNet) of pixels in each image. Our goal is to identify sparse pixel-subsets that contain feature patterns the model identifies as strong class-evidence as it classifies an image. We identify pixels to retain based on sorting by SIS BackSelect [4] (CIFAR-10) or our Batched Gradient BackSelect procedure (ImageNet). These backward selection (BS) pixel-subset images contain the final pixels

	airplane	automobile	bird	cat	deer	dog	frog	horse	ship	truck
airplane										
ResNet18										
ResNet20										
VGG16										
Adv. Robust										

Figure shows masked pixel of each image the model was still able to classify these masked image with >99% confidence.

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.

(with their same RGB values as in the original images) while all other pixels’ values are replaced with zero. Note that we apply backward selection to the function giving the confidence of the *predicted* class from the original model to prevent adding information about the true class for misclassified images, and we use the true labels for training/evaluating models on pixel-subsets. As backward selection is applied locally on each image, the specific pixels retained differ across images.

We train new classifiers on solely these pixel-subsets of training images and evaluate accuracy on corresponding pixel-subsets of test images to determine whether such pixel-subsets are statistically valid for generalization in the benchmark. We use the same training setup and hyperparameters (Section 3.1) without data augmentation of training images (results with data augmentation in Table S1). We consider a model to overinterpret its input when these signals can generalize to test data but lack semantic meaning (Section 3.4).

Author used pixel subset (SIS) both in training and testing set.

### 3.4 Human Classification Benchmark

To evaluate whether sparse pixel-subsets of images can be accurately classified by humans, we asked four participants to classify images containing various degrees of masking. We randomly sampled 100 images from the CIFAR-10 test set (10 images per class) that were correctly and confidently ( $\geq 99\%$  confidence) classified by our models, and for each image, kept only 5%, 30%, or 50% of pixels as ranked by backward selection (all other pixels masked). Backward selection image subsets are sampled across our three models. Since larger subsets of pixels are by construction supersets of smaller subsets identified by the same model, we presented each batch of 100 images in order of increasing subset size and shuffled the order of images within each batch. Users were asked to classify each of the 300 images as one of the 10 classes in CIFAR-10 and were not provided training images. The same task was given to each user (and is shown in Section S5).

## 4 Results

### 4.1 CNNs Classify Images Using Spurious Features

**CIFAR-10.** Figure 1 shows example SIS subsets (threshold 0.99) from CIFAR-10 test images (additional examples in Section S3). These SIS subset images are confidently and correctly classified by each model with  $\geq 99\%$  confidence toward the predicted class. We observe these SIS subsets are highly sparse and the average SIS size at this threshold is  $< 5\%$  of each image (Figure 2), suggesting these CNNs confidently classify images that appear nonsensical to humans (Section 4.3), leading to

The author found out that the classifier were able to predict true class with high confidence even when the SIS input seems non-sensical to humans

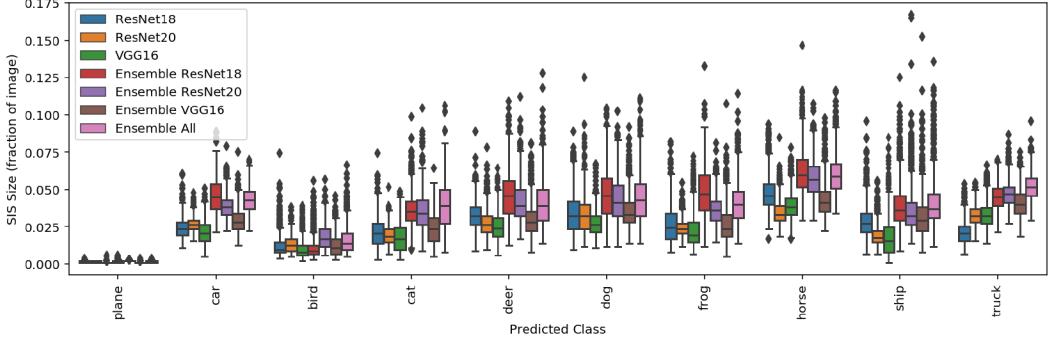


Figure 2: Distribution of SIS size per predicted class by CIFAR-10 models computed on all CIFAR-10 test set images classified with  $\geq 99\%$  confidence (SIS confidence threshold 0.99).

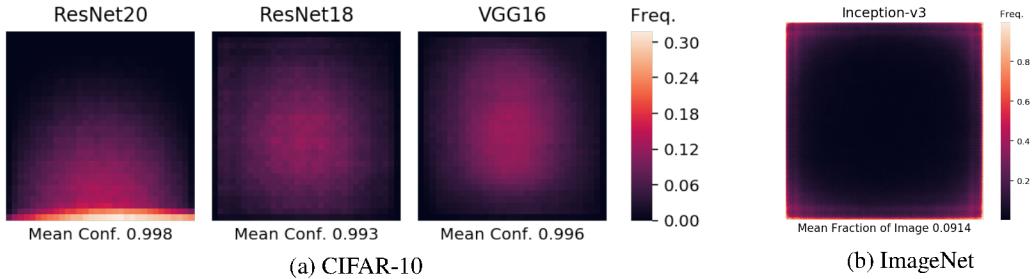


Figure 3: Heatmaps of pixel locations comprising pixel-subsets. Frequency indicates fraction of subsets containing each pixel. **(a)** 5% pixel-subsets across CIFAR-10 test set for each model. Mean confidence indicates confidence on 5% pixel-subsets. **(b)** Sufficient input subsets (confidence threshold 0.9) across ImageNet validation images from Inception v3.

concern about their robustness and generalizability. We also find that SIS size can differ significantly by predicted class (Figure 2).

We retain 5% of pixels in each image using local backward selection and mask the remaining 95% with zeros (Section 3.3) and find models trained on full images classify these pixel-subsets as accurately as full images (Table 1). Figure 3a shows the pixel locations and confidence of these 5% pixel-subsets across all CIFAR-10 test images. We found the concentration of pixels on the bottom border for ResNet20 is a result of tie-breaking during SIS backward selection (Section S4). Moreover, the CNNs are more confident on these pixels subsets than on full images: the mean drop in confidence for the predicted class between original images and these 5% subsets is  $-0.035$  (std dev. = 0.107),  $-0.016$  (0.094), and  $-0.012$  (0.074) computed over all CIFAR-10 test images for our ResNet20, ResNet18, and VGG16 models, respectively, suggesting severe overinterpretation (negative values imply greater confidence on the 5% subsets). We find pixel-subsets chosen via backward selection are significantly more predictive than equally large pixel-subsets chosen uniformly at random from each image (Table 1).

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

While recent work has suggested semantics can be better captured by models that are robust to adversarial inputs that fool standard neural networks via human-imperceptible modifications to images [23, 38], we explore a wide residual network that is adversarially robust for CIFAR-10

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, accuracy is reported 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. Additional results on CIFAR-10.1 are presented in Table S4.

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	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	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
VGG16	5% BS Subsets	5% BS Subsets	94.96 $\pm$ 0.04	75.25 $\pm$ 0.05
		5% Random	51.27 $\pm$ 0.82	45.24 $\pm$ 0.45
		Input Dropout (Full)	94.15 $\pm$ 0.26	80.35 $\pm$ 0.39
	Ensemble (ResNet18)	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

classification [23] and find evidence of overinterpretation (Figure 1). This finding suggests adversarial robustness alone does not prevent models from overinterpreting spurious signals in CIFAR-10.

We also ran Batched Gradient SIS on CIFAR-10 and found edge-heavy sufficient input subsets for CIFAR-10 (Section S4). These heatmap differences are a result of the different valid equivalent sufficient input subsets found by the two SIS discovery algorithms. However, since all sufficient input subsets are validated with a model and guaranteed to be sufficient for classification at the specified threshold, the heatmaps are accurate depictions of what is sufficient for the model to classify images at the threshold. Overinterpretation is independent of the SIS algorithm used because both algorithms produce human-uninterpretable sufficient subsets as shown in the examples.

**ImageNet.** We find models trained on ImageNet images suffer from severe overinterpretation. Figure 4 shows example SIS subsets (threshold 0.9) found via Batched Gradient SIS on images confidently classified by the pre-trained Inception v3 (additional examples in Figures S12–S14). These SIS subsets appear visually nonsensical, yet the network classifies them with  $\geq 90\%$  confidence. We find SIS pixels are concentrated outside of the actual object that determines the class label. For example, in the “pizza” image, the SIS is concentrated on the shape of the plate and the background table, rather than the pizza itself, suggesting the model could generalize poorly on images containing different circular items on a table. In the “giant panda” image, the SIS contains bamboo, which likely appeared in the collection of ImageNet photos for this class. In the “traffic light” and “street sign” images, the SIS consists of pixels in sky, suggesting that autonomous vehicle systems that may depend on these models should be carefully evaluated for overinterpretation pathologies.

Figure 3b shows SIS pixel locations from a random sample of 1000 ImageNet validation images. We find concentration along image borders, suggesting the model relies heavily on image backgrounds and suffers from severe overinterpretation. This is a serious problem as objects determining ImageNet



Figure 4: 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.

classes are often located near image centers, and thus this network fails to focus on salient features. We found the mean fraction of an image required for classification with  $\geq 90\%$  confidence is only 0.0914, and mean SIS size differs significantly by predicted class (Figure S16).

#### 4.2 Sparse Subsets are Real Statistical Patterns

The overconfidence of CNNs for image classification [39] may lead one to wonder whether the observed overconfidence on semantically meaningless SIS subsets is an artifact of calibration rather than true statistical signals in the dataset. We train models on 5% pixel-subsets of CIFAR-10 training images found via backward selection (Section 3.3). We find models trained solely on these pixel-subsets can classify corresponding test image pixel-subsets with minimal accuracy loss compared to models trained on full images (Table 1), and thus these 5% pixel-subsets are valid statistical signals in training images that generalize to the test distribution. As a baseline to the 5% pixel-subsets identified by backward selection, we create variants of all images where the 5% pixel-subsets are selected at random from each image (rather than by backward selection) and use the same random pixel-subsets for training each new model. Models trained on random subsets have significantly lower test accuracy compared to models trained on 5% pixel-subsets from backward selection (Table 1). We observe, however, that random 5% subsets of images still capture enough signal to predict roughly 5 times better than blind guessing, but do not capture nearly enough information for models to make accurate predictions.

We found that the 5% backward selection pixel-subsets did not contain model-specific features, and thus reflected valid predictive signals regardless of the model architecture employed for subset discovery. Our hypothesis was that 5% pixel-subsets discovered with one architecture would provide robust performance when used to train and evaluate a second architecture. We found this hypothesis supported for all six pairs of subset discovery and train-test architectures evaluated (Table S2). These results demonstrate that the highly sparse subsets found via backward selection offer a valid predictive signal in the CIFAR-10 benchmark exploited by models to attain high test accuracy.

We observe similar results on ImageNet. Inception v3 trained on 10% pixel-subsets of ImageNet training images achieves 71.4% top-1 accuracy (mean over 5 runs) on the corresponding pixel-subset ImageNet validation set (Table S7). Additional ImageNet results for Inception v3 and ResNet50, including training and evaluation on random pixel-subsets and pixel-subsets of different architectures, are provided in Table S7.

#### 4.3 Humans Struggle to Classify Sparse Subsets

We find a strong correlation between the fraction of unmasked pixels in each image and human classification accuracy ( $R^2 = 0.94$ , Figure S11). Human accuracy on 5% pixel-subsets of CIFAR-10 images (mean = 19.2%, std dev = 4.8%, Table S6) is significantly lower than on original, unmasked images (roughly 94% [40]), though greater than random guessing, presumably due to correlations between labels and features such as color (e.g., blue sky suggests airplane, ship, or bird).

However, CNNs (even when trained on full images and achieve accuracy on par with human accuracy on full images) classify these sparse image subsets with very high accuracy (Table 1), indicating

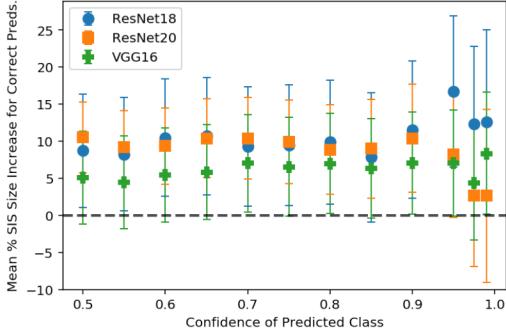


Figure 5: Percentage increase in mean SIS size of correctly classified compared to misclassified CIFAR-10 test images. Positive values indicate larger mean SIS size for correctly classified images. Error bars indicate 95% confidence interval for the difference in means.

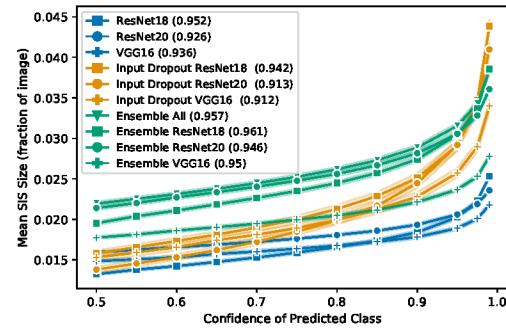


Figure 6: Mean SIS size on CIFAR-10 test images as SIS threshold varies. SIS size indicates fraction of pixels necessary for model to make the same prediction at each confidence threshold. Model accuracies are shown in the legend. 95% confidence intervals are shaded around each mean.

benchmark images contain statistical signals that are not salient to humans. Models solely trained to minimize prediction error may thus latch onto these signals while still accurately generalizing to test data, but may behave counterintuitively when fed images from a different source that does not share these exact statistics. The strong correlation between the size of CIFAR-10 pixel-subsets and the corresponding human classification accuracy suggests larger subsets contain more semantically salient content. Thus, a model whose decisions have larger corresponding SIS subsets presumably exhibits less overinterpretation than one with smaller SIS subsets, as we investigate in Section 4.4.

*Models having larger SIS subset will contain semantically better patterns which can be easily understood by the humans*

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

Given that smaller SIS contain fewer salient features according to human classifiers, models that justify their classifications based on sparse SIS subsets may be limited in terms of attainable accuracy, particularly in out-of-distribution settings. Here, we investigate the relationship between a single model’s predictive accuracy and the size of the SIS subsets in which it identifies class-evidence. We draw no conclusions between models as they are uncalibrated (additional results of SIS from calibrated models are presented in Section S4). For each of our three classifiers, we compute the average SIS size increase for correctly classified images as compared to incorrectly classified images (expressed as a percentage). We find SIS subsets of correctly classified images are consistently significantly larger than those of misclassified images at all SIS confidence thresholds for both CIFAR-10 test images (Figure 5) and CIFAR-10-C OOD images (Figure S3). This is especially striking given model confidence is uniformly lower on the misclassified inputs (Figure S4). Lower confidence would normally imply a larger SIS subset at a given confidence level, as one expects fewer pixels can be masked before the model’s confidence drops below the SIS threshold. Thus, we can rule out overall model confidence as an explanation of the smaller SIS of misclassified images. This result suggests the sparse SIS subsets highlighted in this paper are not just a curiosity, but may be leading to poor generalization on real images.

#### 4.5 Mitigating Overinterpretation

**Ensembling.** Model ensembling is known to improve classification performance [41, 42]. As we found pixel-subset size to be strongly correlated with human pixel-subset classification accuracy (Section 4.3), our metric for measuring how much ensembling may alleviate overinterpretation is the increase in SIS subset size. We find ensembling uniformly increases test accuracy as expected but also increases the SIS size (Figure 6), hence mitigating overinterpretation.

We conjecture the cause of both the increase in the accuracy and SIS size for ensembles is the same. We observe that SIS subsets are generally not transferable from one model to another — i.e., an SIS for one model is rarely an SIS for another (Section 4.1). Thus, different models rely on different independent signals to arrive at the same prediction. An ensemble bases its prediction on multiple

*One of the techniques mitigate overinterpretation is ensembling as every model has separate set of SIS, because of which the overall accuracy increase as it's prediction are dependent on variety of regions in image.*

such signals, increasing predictive accuracy and SIS subset size by requiring simultaneous activation of multiple independently trained feature detectors. We find SIS subsets of the ensemble are larger than the SIS of its individual members (examples in Figure S2).

**Input Dropout.** We apply input dropout [43] to both train and test images. We retain each input pixel with probability  $p = 0.8$  and set the values of dropped pixels to zero. We find a small decrease in CIFAR-10 test accuracy for models regularized with input dropout though find a significant ( $\sim 6\%$ ) increase in OOD test accuracy on CIFAR-10-C images (Table 1, Figure S5). Figure 6 shows a corresponding increase in SIS subset size for these models, suggesting input dropout applied at train and test time helps to mitigate overinterpretation. We conjecture that random dropout of input pixels disrupts spurious signals that lead to overinterpretation.

## 5 Discussion

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

Our results call into question model interpretability methods whose outputs are encouraged to align with prior human beliefs of proper classifier operating behavior [44]. Given the existence of non-salient pixel-subsets that alone suffice for correct classification, a model may solely rely on such patterns. In this case, an interpretability method that faithfully describes the model should output these nonsensical rationales, whereas interpretability methods that bias rationales toward human priors may produce results that mislead users to think their models behave as intended.

Mitigating overinterpretation and the broader task of ensuring classifiers are accurate for the right reasons remain significant challenges for ML. While we identify strategies for partially mitigating overinterpretation, additional research needs to develop ML methods that rely exclusively on well-formed interpretable inputs, and methods for creating training data that do not contain spurious signals. One alternative is to regularize CNNs by constraining the pixel attributions generated via a saliency map [45–47]. Unfortunately, such methods require a human annotator to highlight the correct pixels as an auxiliary supervision signal. Saliency maps have also been shown to provide unreliable insights into model operating behavior and must be interpreted as approximations [48]. In contrast, our SIS subsets constitute actual pathological examples that have been misconstrued by the model. An important application of our methods is the evaluation of training datasets to ensure decisions are made on interpretable rather than spurious signals. We found popular image datasets contain such spurious signals, and the resulting overinterpretation may be difficult to overcome with ML methods alone.

## Acknowledgments and Disclosure of Funding

This work was supported by Schmidt Futures and the National Institutes of Health [R01CA218094].

## Author Contributions

All authors contributed to conceptualization, methodology, formal analysis, and writing. BC led execution of the experiments.

## References

- [1] Neel V. Patel. *Why Doctors Aren't Afraid of Better, More Efficient AI Diagnosing Cancer*, 2017. URL <https://www.thedailybeast.com/why-doctors-arent-afraid-of-better-more-efficient-ai-diagnosing-cancer>. Accessed September 27, 2020.
- [2] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [3] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- [4] Brandon Carter, Jonas Mueller, Siddhartha Jain, and David Gifford. What made you do this? Understanding black-box decisions with sufficient input subsets. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 567–576, 2019.
- [5] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. In *Advances in Neural Information Processing Systems*, 2019.
- [6] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1521–1528. IEEE, 2011.
- [7] Tatiana Tommasi, Novi Patricia, Barbara Caputo, and Tinne Tuytelaars. A deeper look at dataset bias. In *Domain Adaptation in Computer Vision Applications*, pages 37–55. Springer, 2017.
- [8] Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal perturbations and smooth masks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2950–2958, 2019.
- [9] Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE Transactions on Neural Networks and Learning Systems*, 28(11):2660–2673, 2016.
- [10] Chirag Agarwal and Anh Nguyen. Explaining image classifiers by removing input features using generative models. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [11] Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and Payel Das. Explanations based on the missing: towards contrastive explanations with pertinent negatives. In *Advances in Neural Information Processing Systems*, 2018.
- [12] Amir Rosenfeld, Richard Zemel, and John K Tsotsos. The elephant in the room. *arXiv preprint arXiv:1808.03305*, 2018.
- [13] Rakshith Shetty, Bernt Schiele, and Mario Fritz. Not using the car to see the sidewalk—quantifying and controlling the effects of context in classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8218–8226, 2019.
- [14] Krishna Kumar Singh, Dhruv Mahajan, Kristen Grauman, Yong Jae Lee, Matt Feiszli, and Deepti Ghadiyaram. Don't judge an object by its context: Learning to overcome contextual bias. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11070–11078, 2020.
- [15] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really learn. *Nature Communications*, 10(1):1–8, 2019.
- [16] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Texture and art with deep neural networks. *Current Opinion in Neurobiology*, 46:178–186, 2017.
- [17] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019.
- [18] Wieland Brendel and Matthias Bethge. Approximating CNNs with Bag-of-local-Features models works surprisingly well on ImageNet. In *International Conference on Learning Representations*, 2019.

- [19] Robert Geirhos, Carlos R Medina Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A Wichmann. Generalisation in humans and deep neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- [20] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 427–436, 2015.
- [21] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- [22] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 427–436, 2015.
- [23] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- [24] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems*, 2019.
- [25] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems*, 2019.
- [26] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- [27] Shi Feng, Eric Wallace, Alvin Grissom II, Mohit Iyyer, Pedro Rodriguez, and Jordan Boyd-Graber. Pathologies of neural models make interpretations difficult. In *Empirical Methods in Natural Language Processing*, 2018.
- [28] Timothy Niven and Hung-Yu Kao. Probing neural network comprehension of natural language arguments. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4658–4664, 2019.
- [29] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2019.
- [30] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do CIFAR-10 classifiers generalize to CIFAR-10? *arXiv preprint arXiv:1806.00451*, 2018.
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European Conference on Computer Vision*. Springer, 2016.
- [33] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [34] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International Conference on Machine Learning*, 2013.
- [35] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped dqn. In *Advances in Neural Information Processing Systems*, 2016.
- [36] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016.
- [37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, 2019.
- [38] Shibani Santurkar, Andrew Ilyas, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Image synthesis with a single (robust) classifier. In *Advances in Neural Information Processing Systems*, 2019.

- [39] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, 2017.
- [40] Andrej Karpathy. Lessons learned from manually classifying CIFAR-10. *Published online at http://karpathy.github.io/2011/04/27/manually-classifying-cifar10*, 2011.
- [41] King-Shy Goh, Edward Chang, and Kwang-Ting Cheng. SVM binary classifier ensembles for image classification. In *Proceedings of the Tenth International Conference on Information and Knowledge Management*, pages 395–402. ACM, 2001.
- [42] Cheng Ju, Aurélien Bibaut, and Mark van der Laan. The relative performance of ensemble methods with deep convolutional neural networks for image classification. *Journal of Applied Statistics*, 45(15):2800–2818, 2018.
- [43] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [44] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, 2018.
- [45] Andrew Slavin Ross, Michael C Hughes, and Finale Doshi-Velez. Right for the right reasons: training differentiable models by constraining their explanations. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 2662–2670, 2017.
- [46] Becks Simpson, Francis Dutil, Yoshua Bengio, and Joseph Paul Cohen. GradMask: Reduce overfitting by regularizing saliency. *arXiv preprint arXiv:1904.07478*, 2019.
- [47] Joseph D Viviano, Becks Simpson, Francis Dutil, Yoshua Bengio, and Joseph Paul Cohen. Saliency is a possible red herring when diagnosing poor generalization. In *International Conference on Learning Representations*, 2021.
- [48] Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. The (un)reliability of saliency methods. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 267–280. Springer, 2019.
- [49] Antonio Torralba, Rob Fergus, and William T Freeman. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(11):1958–1970, 2008.
- [50] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4700–4708, 2017.
- [51] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [52] Muzammal Naseer, Kanchana Ranasinghe, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers. *arXiv preprint arXiv:2105.10497*, 2021.

---

# Supplementary Material: Overinterpretation reveals image classification model pathologies

---

**Brandon Carter**  
MIT CSAIL  
bcarter@csail.mit.edu

**Siddhartha Jain**  
MIT CSAIL

**Jonas Mueller**  
Amazon Web Services

**David Gifford**  
MIT CSAIL  
gifford@mit.edu

## Contents

<b>S1 Details of Batched Gradient SIS Algorithm</b>	<b>2</b>
<b>S2 Model Implementation and Training Details</b>	<b>4</b>
<b>S3 Additional Examples of CIFAR-10 Sufficient Input Subsets</b>	<b>5</b>
S3.1 SIS of Individual Networks . . . . .	5
S3.2 Ensemble Sufficient Input Subsets . . . . .	6
<b>S4 Additional Results on CIFAR-10</b>	<b>7</b>
S4.1 Training on Pixel-Subsets With Data Augmentation . . . . .	7
S4.2 Training on Pixel-Subsets With Different Architectures . . . . .	7
S4.3 Additional Results for Models Trained on Pixel-Subsets . . . . .	7
S4.4 Additional Results for SIS Size and Model Accuracy . . . . .	9
S4.5 Additional Results for Input Dropout . . . . .	10
S4.6 Results on CIFAR-10.1 . . . . .	11
S4.7 SIS and Calibrated Models . . . . .	11
S4.8 SIS with Random Tie-breaking . . . . .	12
S4.9 Confidence Curves for SIS Backward Selection on CIFAR-10 . . . . .	13
S4.10Batched Gradient SIS on CIFAR-10 . . . . .	14
<b>S5 Details of Human Classification Benchmark</b>	<b>15</b>
<b>S6 Additional Results of ImageNet Overinterpretation</b>	<b>17</b>
S6.1 Training CNNs on ImageNet Pixel-Subsets . . . . .	17
S6.2 Additional Examples of SIS on ImageNet . . . . .	17
S6.3 SIS Size by Class . . . . .	22
S6.4 SIS for Vision Transformers . . . . .	22
<b>S7 NeurIPS Paper Checklist</b>	<b>23</b>

## S1 Details of Batched Gradient SIS Algorithm

It is computationally infeasible to scale the original backward selection procedure of SIS [4] to ImageNet. As each ImageNet image contains  $299 \times 299 = 89401$  pixels, running backward selection to find one SIS for an image would require  $\sim 4$  billion forward passes through the network. Here we introduce a more efficient gradient-based approximation to the original SIS procedure (via **Batched Gradient SIScollection**, **Batched Gradient BackSelect**, and **Batched Gradient FindSIS**) that allows us to find SIS on larger ImageNet images in a reasonable time. The **Batched Gradient SIScollection** procedure described below identifies a complete collection of disjoint masks for an input  $\mathbf{x}$ , where each mask  $M$  specifies a pixel-subset of the input  $\mathbf{x}_S = \mathbf{x} \odot (1 - M)$  such that  $f(\mathbf{x}_S \geq \tau)$ . Here  $f$  outputs the probability assigned by the network to its predicted class (i.e., its confidence).

The idea behind our approximation algorithm is two-fold: (1) Instead of separately masking every remaining pixel to find the least critical pixel (whose masking least reduces the confidence in the network’s prediction), we use the *gradient* with respect to the mask as a means of ordering. (2) Instead of masking just 1 pixel per iteration, we mask larger subsets of  $k \geq 1$  pixels per iteration. More formally, let  $\mathbf{x}$  be an image of dimensions  $H \times W \times C$  where  $H$  is the height,  $W$  the width, and  $C$  the channel. Let  $f(\mathbf{x})$  be the network’s confidence on image  $\mathbf{x}$  and  $\tau$  the target SIS confidence threshold. Recall that we only compute SIS for images where  $f(\mathbf{x}) \geq \tau$ . Let  $M$  be the mask with dimensions  $H \times W$  with 0 indicating an unmasked feature (pixel) and 1 indicating a masked feature. We initialize  $M$  as all 0s (all features unmasked). At iteration  $i$ , we compute the gradient of  $f$  with respect to the input pixels and mask  $\nabla M = \nabla_M f(\mathbf{x} \odot (1 - M))$ . Here  $M$  is the current mask updated after each iteration. In each iteration, we find the block of  $k$  features to mask,  $G^*$ , chosen in descending order by value of entries in  $\nabla M$ . The mask is updated after each iteration by masking this block of  $k$  features until all features have been masked. Given  $p$  input features, our **Batched Gradient SIScollection** procedure returns  $j$  sufficient input subsets in  $\mathcal{O}(\frac{p}{k} \cdot j)$  evaluations of  $\nabla f$  (as opposed to  $\mathcal{O}(p^2 j)$  evaluations of  $f$  in the original SIS procedure [4]).

We use  $k = 100$  in this paper, which allows us to find one SIS for each of 32 ImageNet images (i.e., a mini-batch) in  $\sim 1\text{-}2$  minutes using **Batched Gradient FindSIS**. Note that while our algorithm is an approximate procedure, the pixel-subsets produced are real sufficient input subsets, i.e., they always satisfy  $f(\mathbf{x}_S \geq \tau)$ . For CIFAR-10 images (which are smaller in size), we use the original SIS procedure from [4]. For both datasets, we treat all channels of each pixel as a single feature.

---

**Algorithm 1: Batched Gradient SIScollection**

---

**Input:** function  $f$ , input  $\mathbf{x}$ , threshold  $\tau$ , batch size  $k$  (number of pixels)  
 $M = \mathbf{0}$   
**for**  $j = 1, 2, \dots$  **do**  
     $R = \text{Batched Gradient BackSelect}(f, \mathbf{x}, M, k)$   
     $M_j = \text{Batched Gradient FindSIS}(f, \mathbf{x}, \tau, R)$   
     $M \leftarrow M + M_j$   
    **if**  $f(\mathbf{x} \odot (1 - M)) < \tau$  **then**  
        **return**  $M_1, \dots, M_{j-1}$   
    **end if**  
**end for**

---

**Algorithm 2: Batched Gradient BackSelect**

---

**Input:** function  $f$ , input  $\mathbf{x}$ , mask  $M$ , batch size  $k$  (number of pixels)  
 $R = \text{empty stack}$   
**while**  $M \neq \mathbf{1}$  **do**  
     $G^* = \text{Top}_k(\nabla_M f(\mathbf{x} \odot (1 - M)))$   
    Update  $M \leftarrow M + G^*$   
    Push  $G^*$  onto top of  $R$   
**end while**  
**return**  $R$

---

**Algorithm 3: Batched Gradient FindSIS**

---

**Input:** function  $f$ , input  $\mathbf{x}$ , threshold  $\tau$ , stack  $R$   
 $M = \mathbf{1}$   
**while**  $f(\mathbf{x} \odot (1 - M)) < \tau$  **do**  
    Pop  $G$  from top of  $R$   
    Update  $M \leftarrow M - G$   
**end while**  
**if**  $f(\mathbf{x} \odot (1 - M)) \geq \tau$  **then**  
    **return**  $M$   
**else**  
    **return** *None*  
**end if**

---

## S2 Model Implementation and Training Details

### CIFAR-10 Models

We first describe the implementation and training details for the CIFAR-10 models used in this paper (Section 3.1). The ResNet20 architecture [31] has 16 initial filters and a total of 0.27M parameters. ResNet18 [32] has 64 initial filters and contains 11.2M parameters. The VGG16 architecture [33] uses batch normalization and contains 14.7M parameters.

All models are trained for 200 epochs with a batch size of 128. We minimize cross-entropy via SGD with Nesterov momentum [34] using momentum of 0.9 and weight decay of 5e-4. The learning rate is initialized as 0.1 and is reduced by a factor of 5 after epochs 60, 120, and 160. Datasets are normalized using per-channel mean and standard deviation, and we use standard data augmentation strategies consisting of random crops and horizontal flips [32].

The adversarially robust model we evaluated is the `adv_trained` model of Madry et al. [23], available on GitHub<sup>1</sup>.

To apply the SIS procedure to CIFAR-10 images, we use an implementation available on GitHub<sup>2</sup>. For confidently classified images on which we run SIS, we find one sufficient input subset per image using the FindSIS procedure. When masking pixels, we mask all channels of each pixel as a single feature.

### ImageNet Models

For finding SIS, we use pre-trained models (Inception v3 [36] and ResNet50 [31]) provided by PyTorch [37] in the `torchvision` package (PyTorch version 1.4.0, `torchvision` version 0.5.0).

When training new ImageNet classifiers, we adopt model implementations and training scripts from PyTorch [37], obtained from GitHub<sup>3</sup>. Models are trained for 90 epochs using batch size 256 (Inception-v3) or 512 (ResNet50). We minimize cross-entropy via SGD using momentum of 0.9 and weight decay of 1e-4. The learning rate is initialized as 0.1 and reduced by a factor of 10 every 30 epochs. Datasets are normalized using per-channel mean and standard deviation. For Inception v3, images are cropped to 299 x 299 pixels. For ResNet50, images are cropped to 224 x 224. When training Inception v3, we define the model using the `aux_logits=False` argument. We do not use data augmentation when training models on pixel-subsets of images.

### Hardware Details

Each CIFAR-10 model is trained on 1 NVIDIA GeForce RTX 2080 Ti GPU. Once models are trained, SIS are computed across multiple GPUs (by parallelizing over individual images). Each SIS (for 1 CIFAR-10 image) takes roughly 30-60 seconds to compute (depending on the model architecture).

ImageNet models are trained on 2–3 NVIDIA Titan RTX GPUs. For finding SIS from pre-trained ImageNet models, we run Batched Gradient BackSelect for batches of 32 images across 10 NVIDIA GeForce RTX 2080 Ti GPUs, which takes roughly 1-2 minutes per batch (details in Section S1).

---

<sup>1</sup>[https://github.com/MadryLab/cifar10\\_challenge](https://github.com/MadryLab/cifar10_challenge)

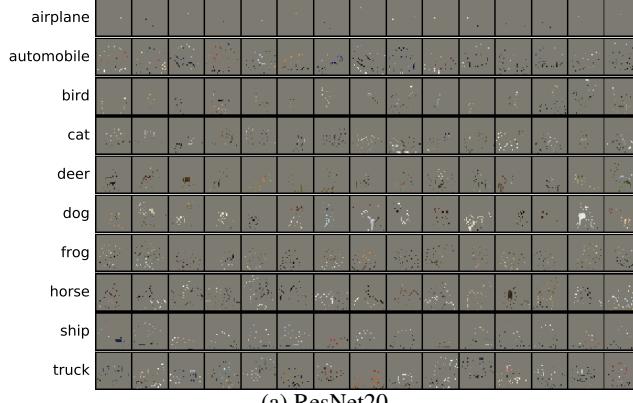
<sup>2</sup>[https://github.com/google-research/google-research/blob/master/sufficient\\_input\\_subsets/sis.py](https://github.com/google-research/google-research/blob/master/sufficient_input_subsets/sis.py)

<sup>3</sup><https://github.com/pytorch/examples/blob/master/imagenet/main.py>

### S3 Additional Examples of CIFAR-10 Sufficient Input Subsets

#### S3.1 SIS of Individual Networks

Figure S1 shows a sample of SIS for each of our three architectures. These images were randomly sampled among all CIFAR-10 test images confidently (confidence  $\geq 0.99$ ) predicted to belong to the class written on the left. Out of 10000 CIFAR-10 test images, 8596 were predicted with  $\geq 99\%$  confidence by ResNet18 (7829 by ResNet20, 9048 by VGG16). SIS are computed under a threshold of 0.99, so all images shown in this figure are classified with probability  $\geq 99\%$  confidence as belonging to the listed class.



(a) ResNet20



(b) ResNet18



(c) VGG16

Figure S1: Examples of SIS (threshold 0.99) on random sample of CIFAR-10 test images (15 per class, different random sample for each architecture). All images shown here are predicted to belong to the listed class with  $\geq 99\%$  confidence.

### S3.2 Ensemble Sufficient Input Subsets

Figure S2 shows examples of SIS from one of our model ensembles (a homogeneous ensemble of ResNet18 networks, see Section 3.1), along with corresponding SIS for the same image from each of the five member networks in the ensemble. We use a SIS threshold of 0.99, so all images are classified with  $\geq 99\%$  confidence. These examples highlight how the ensemble SIS are larger and draw class-evidence from the individual members' SIS.

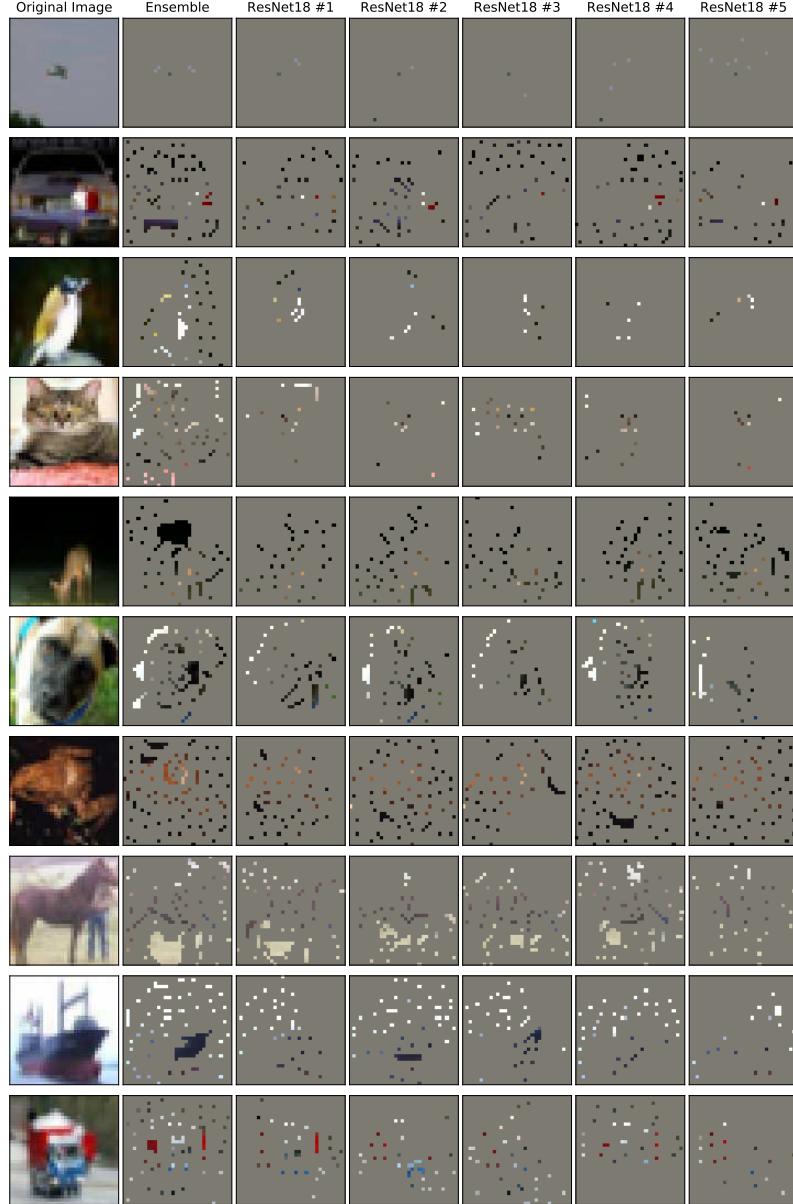


Figure S2: Examples of SIS (threshold 0.99) from the ResNet18 homogeneous ensemble (Section 3.1) and its member models. Each row shows original CIFAR-10 image (left), followed by SIS from the ensemble (second column) and the SIS from each of its 5 member networks (remaining columns). Each image shown is classified with  $\geq 99\%$  confidence by its respective network.

## S4 Additional Results on CIFAR-10

### S4.1 Training on Pixel-Subsets With Data Augmentation

Table S1 presents results similar to those in Section 4.2 and Table 1, but where models are trained on 5% pixel-subsets with data augmentation (as described in Section S2). We find training without data augmentation slightly improves accuracy when training classifiers on 5% pixel-subsets of CIFAR-10.

Table S1: Accuracy of CIFAR-10 classifiers trained and evaluated on full images, 5% backward selection (BS) pixel-subsets, and 5% random pixel-subsets *with* data augmentation (+). Accuracy is reported as mean  $\pm$  standard deviation (%) over five runs.

Model	Train On	Evaluate On	CIFAR-10 Test Acc.	CIFAR-10-C Acc.
ResNet20	5% BS Subsets (+)	5% BS Subsets	92.26 $\pm$ 0.01	70.21 $\pm$ 0.14
	5% Random (+)	5% Random	48.87 $\pm$ 0.41	42.66 $\pm$ 0.15
ResNet18	5% BS Subsets (+)	5% BS Subsets	94.51 $\pm$ 0.38	74.91 $\pm$ 0.41
	5% Random (+)	5% Random	49.03 $\pm$ 0.92	42.97 $\pm$ 0.82
VGG16	5% BS Subsets (+)	5% BS Subsets	91.17 $\pm$ 0.04	71.82 $\pm$ 0.13
	5% Random (+)	5% Random	51.32 $\pm$ 1.35	44.56 $\pm$ 0.96

### S4.2 Training on Pixel-Subsets With Different Architectures

Table S2 presents results of training and evaluating models on 5% pixel-subsets drawn from different architectures. Models were trained without data augmentation on subsets from one replicate of each base architecture. We find accuracy from training and evaluating a model on 5% pixel-subsets of images derived from a different architecture is commensurate with accuracy of training and evaluating a new model of the same type on those subsets (Table 1).

Table S2: Accuracy of CIFAR-10 classifiers trained and evaluated on 5% backward selection (BS) pixel-subsets from different architectures. Accuracy is reported as mean  $\pm$  standard deviation (%) over five runs.

5% Subsets from Model	Model Trained	CIFAR-10 Test Acc.	CIFAR-10-C Acc.
ResNet20	ResNet18	92.53 $\pm$ 0.02	70.56 $\pm$ 0.04
	VGG16	92.47 $\pm$ 0.02	70.42 $\pm$ 0.14
ResNet18	ResNet20	94.88 $\pm$ 0.03	75.14 $\pm$ 0.10
	VGG16	94.88 $\pm$ 0.05	75.13 $\pm$ 0.09
VGG16	ResNet20	92.05 $\pm$ 0.14	73.01 $\pm$ 0.08
	ResNet18	92.57 $\pm$ 0.10	73.33 $\pm$ 0.21

### S4.3 Additional Results for Models Trained on Pixel-Subsets

Table S3 presents results of models trained on 5% backward selection (BS) or random pixel-subsets of CIFAR-10 training images, evaluated on full (original) CIFAR-10 test images. While accuracies are generally significantly higher than random guessing, we note that full images are highly out-of-distribution for a model trained on images with only 5% pixel-subsets and hence such a model cannot properly generalize to full images. Further, the model trained on 5% images may not rely on the same features as the model trained on full images as it is trained on a substantially different training set.

Table S3: Accuracy of CIFAR-10 classifiers trained on 5% backward selection (BS) or random pixel-subsets with (+) and without (−) data augmentation. Accuracy is reported as mean  $\pm$  standard deviation (%) over five runs.

Model	Train On	Evaluate On	CIFAR-10 Test Acc.	CIFAR-10-C Acc.
ResNet20	5% BS Subsets (−)	Full Images	$21.02 \pm 1.57$	$17.50 \pm 1.15$
	5% Random (−)	Full Images	$38.66 \pm 3.31$	$36.40 \pm 2.73$
	5% BS Subsets (+)	Full Images	$10.87 \pm 1.50$	$10.75 \pm 1.32$
	5% Random (+)	Full Images	$37.08 \pm 3.51$	$33.78 \pm 2.81$
ResNet18	5% BS Subsets (−)	Full Images	$20.86 \pm 2.74$	$18.20 \pm 1.43$
	5% Random (−)	Full Images	$26.05 \pm 7.59$	$25.03 \pm 6.41$
	5% BS Subsets (+)	Full Images	$11.83 \pm 1.74$	$11.48 \pm 1.15$
	5% Random (+)	Full Images	$20.98 \pm 4.61$	$20.35 \pm 3.56$
VGG16	5% BS Subsets (−)	Full Images	$41.63 \pm 3.55$	$30.34 \pm 1.97$
	5% Random (−)	Full Images	$25.73 \pm 6.08$	$23.56 \pm 4.39$
	5% BS Subsets (+)	Full Images	$14.32 \pm 3.40$	$13.22 \pm 2.01$
	5% Random (+)	Full Images	$27.58 \pm 3.96$	$24.92 \pm 3.10$

#### S4.4 Additional Results for SIS Size and Model Accuracy

Figure S3 shows percentage increase in mean SIS size for correctly classified images compared to misclassified images from the CIFAR-10-C dataset.

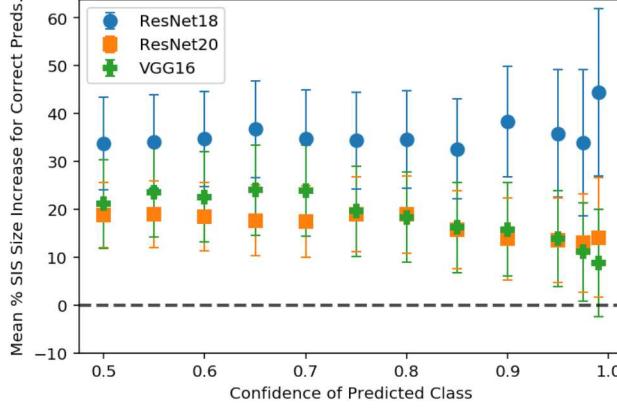


Figure S3: Percentage increase in mean SIS size of correctly classified images compared to misclassified images from a random sample of CIFAR-10-C test set. Positive values indicate larger mean SIS size for correctly classified images. Error bars indicate 95% confidence interval for the difference in means.

Figure S4 shows the mean confidence of each group of correctly and incorrectly classified images that we consider at each confidence threshold (at each confidence threshold along the x-axis, we evaluate SIS size in Figure 5 on the set of images that originally were classified with at least that level of confidence). We find model confidence is uniformly lower on the misclassified inputs.

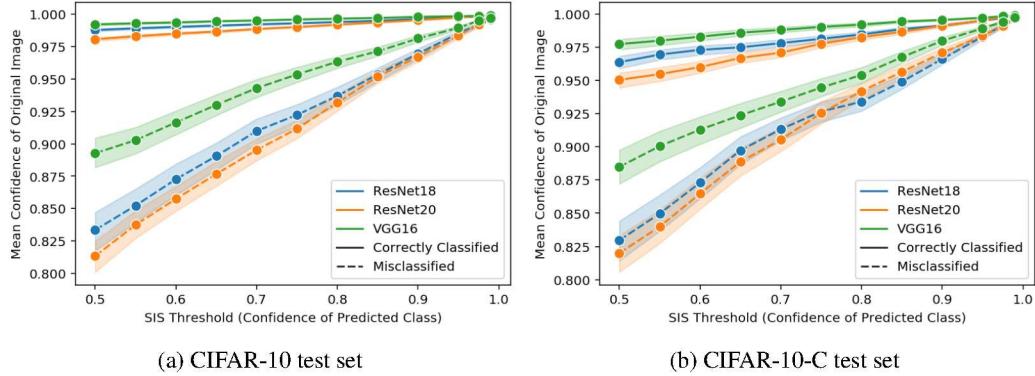


Figure S4: Mean confidence of correctly vs. incorrectly classified images for each corresponding SIS threshold we evaluate in Figure 5 across the (a) CIFAR-10 test set and (b) our random sample of the CIFAR-10-C test set. Shaded region indicates 95% confidence interval.

## S4.5 Additional Results for Input Dropout

Figure S5 shows the accuracy improvement on each individual corruption of the CIFAR-10-C out-of-distribution test set for models trained with input dropout (Section 4.5) compared to original models.

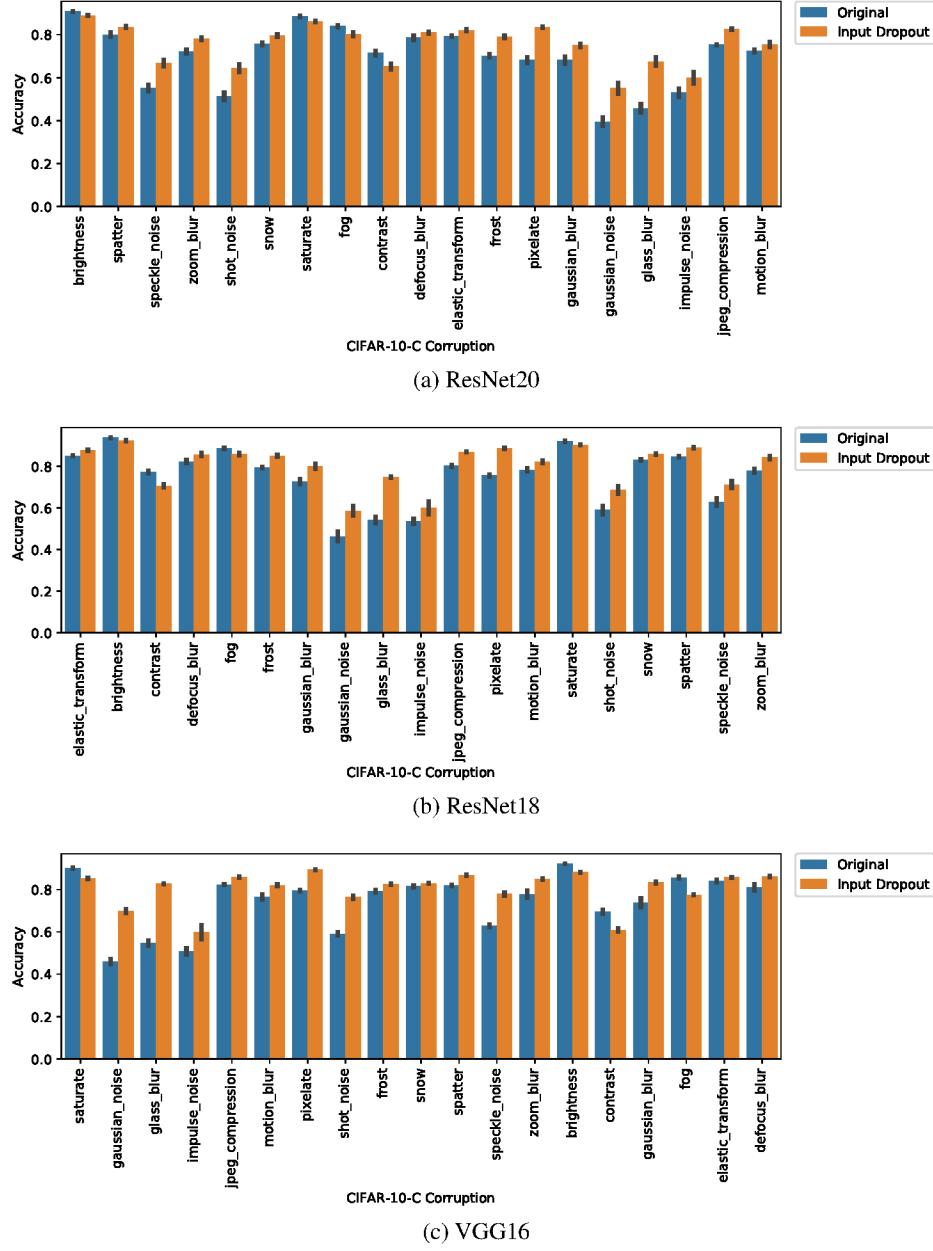


Figure S5: Accuracy on individual corruptions of CIFAR-10-C out-of-distribution images for original models and models trained with input dropout (Section 4.5). Accuracy is given as mean  $\pm$  standard deviation over five replicate models.

#### S4.6 Results on CIFAR-10.1

Table S4 reports accuracy of the models from Section 4.2 computed on the CIFAR-10.1 v6 dataset [30], which contains 2000 class-balanced images drawn from the Tiny Images repository [49] in a similar fashion to that of CIFAR-10, though Recht et al. [30] found a large drop in classification accuracy on these images.

Table S4: Accuracy of CIFAR-10 classifiers trained and evaluated on full images, 5% backward selection (BS) pixel-subsets, and 5% random pixel-subsets reported on CIFAR-10.1 v6 dataset (evaluating models from Section 4.2 that were trained on full images or 5% subsets of the CIFAR-10 train set). Where possible, accuracy is reported 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-10.1 Acc.	
ResNet20	Full Images	Full Images	83.98 $\pm$ 0.68	
		5% BS Subsets	82.80	
		5% Random	10.00 $\pm$ 0.00	
	5% BS Subsets	5% BS Subsets	82.56 $\pm$ 0.07	
		5% Random	39.78 $\pm$ 1.27	
	Input Dropout (Full)	Input Dropout (Full)	81.88 $\pm$ 0.44	
	Full Images	Full Images	88.89 $\pm$ 0.45	
		5% BS Subsets	89.35	
		5% Random	10.06 $\pm$ 0.11	
ResNet18	5% BS Subsets	5% BS Subsets	89.49 $\pm$ 0.04	
		5% Random	39.45 $\pm$ 1.02	
	Input Dropout (Full)	Input Dropout (Full)	86.28 $\pm$ 0.33	
	Full Images	Full Images	86.23 $\pm$ 0.79	
VGG16		5% BS Subsets	86.45	
		5% Random	9.78 $\pm$ 0.26	
5% BS Subsets	5% BS Subsets	85.61 $\pm$ 0.19		
	5% Random	40.98 $\pm$ 1.27		
Input Dropout (Full)	Input Dropout (Full)	81.00 $\pm$ 0.65		
Ensemble (ResNet18)	Full Images	Full Images	90.30	
		5% Random	10.05	

#### S4.7 SIS and Calibrated Models

We calibrated one model of each architecture class after training using Temperature Scaling [39] based on an implementation available on GitHub<sup>4</sup>. The CIFAR-10 test set was randomly split into a 5k validation set (for optimization of the temperature parameter) and a 5k held-out test set (for final evaluation of ECE). Table S5 shows the Expected Calibration Error (ECE) of each model on held-out test images before and after calibration, as well as mean SIS size using confidence threshold 0.99 computed on the entire CIFAR-10 test set. We find that while the mean SIS size (for test images that the re-calibrated model can classify with  $\geq 99\%$  confidence) does increase slightly, the resulting SIS subsets are still semantically meaningless and far below the threshold of SIS size where humans can meaningfully start to classify CIFAR images with any degree of accuracy (Figure S6). We note that one of the key findings of our paper is that even when we compute SIS subsets from uncalibrated models, those subsets still contain enough signal for training entirely new classifiers that can generalize as well to the corresponding test subsets (Section 4.2).

<sup>4</sup>[https://github.com/gpleiss/temperature\\_scaling](https://github.com/gpleiss/temperature_scaling)

Table S5: Results of model calibration by temperature scaling. Expected Calibration Error (ECE) is computed on a held-out set of 5k CIFAR-10 test images. SIS are computed using a threshold of 0.99 on all CIFAR-10 test images classified with  $\geq 99\%$  confidence (and corresponding number of such images listed). SIS size is given as mean  $\pm$  standard deviation.

Model	ECE (%)	SIS Size (% of Image)	Num. Images Pred. $\geq 0.99$
ResNet20 Uncalibrated	3.91	$2.36 \pm 1.21$	7829
ResNet20 Calibrated	0.91	$2.94 \pm 1.39$	5805
ResNet18 Uncalibrated	2.49	$2.53 \pm 1.53$	8596
ResNet18 Calibrated	1.00	$3.54 \pm 1.94$	5934
VGG16 Uncalibrated	4.95	$2.18 \pm 1.37$	9048
VGG16 Calibrated	1.56	$8.26 \pm 2.86$	23



Figure S6: Examples of SIS (threshold 0.99) on sample of CIFAR-10 test images from calibrated models. All images shown are predicted to belong to the listed class with  $\geq 99\%$  confidence.

#### S4.8 SIS with Random Tie-breaking

We suspect the concentration of pixels on the bottom border for ResNet20 (Figure 3a) is a result of tie-breaking during backward selection of the SIS procedure. To explore this hypothesis, we modified the tie-breaking procedure to randomly (rather than deterministically) break ties during SIS backward selection by adding random Gaussian noise ( $\mu = 0$ ,  $\sigma^2 = 1e-12$ ) to the model’s outputs for each remaining masked pixel at each iteration of backward selection. For each image in a sample of 1000 CIFAR-10 test images, we repeated this randomization procedure three times and found the resulting heatmap of 5% backward selection pixel-subsets for ResNet20 more concentrated in the image centers rather than bottom border (Figure S7).

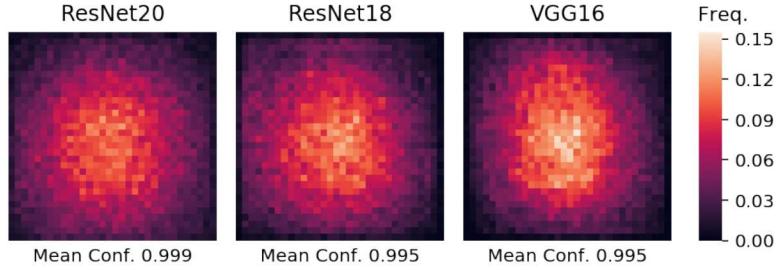


Figure S7: Heatmap of pixel locations comprising 5% backward selection pixel-subsets computed on a set of 1000 CIFAR-10 test set images with random tie-breaking during backward selection.

#### S4.9 Confidence Curves for SIS Backward Selection on CIFAR-10

Figure S8 shows the predicted confidence on the remaining pixels at each step of SIS backward selection for the entire CIFAR-10 test set for each architecture trained on CIFAR-10.

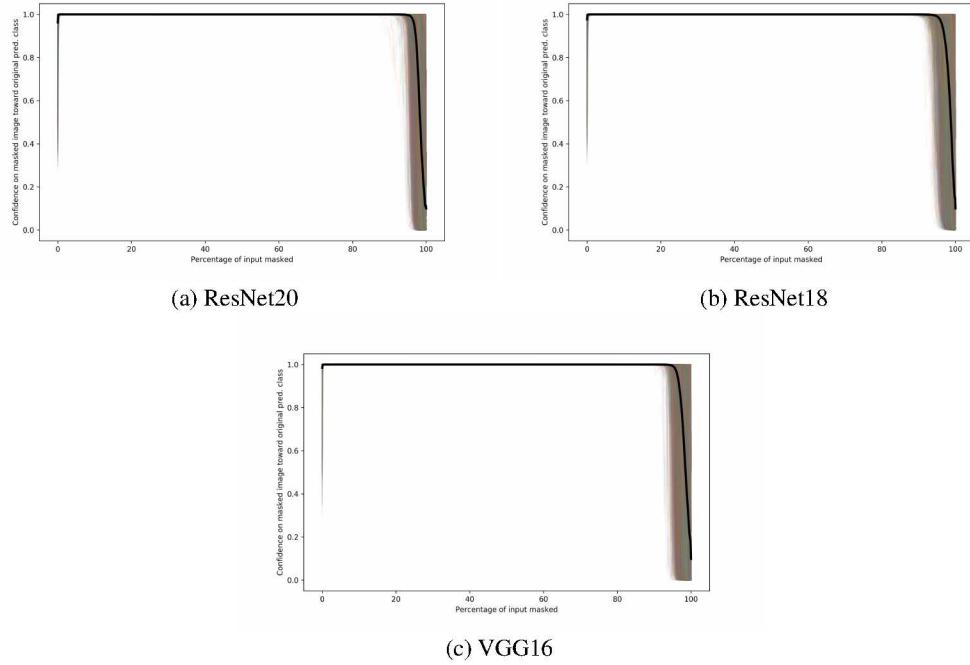


Figure S8: Prediction history on remaining (unmasked) pixels at each step of the SIS backward selection procedure for all CIFAR-10 test set images. Black line depicts mean confidence at each step.

#### S4.10 Batched Gradient SIS on CIFAR-10

We also ran Batched Gradient SIS on the entire CIFAR-10 test set for ResNet18 and found Batched Gradient SIS produced edge-heavy heatmaps for CIFAR-10 (Figure S9a). For CIFAR-10, we set  $k = 1$  to remove a single pixel per iteration of Batched Gradient SIS. These heatmap differences (compared to Figure 3) are a result of the different valid equivalent SIS subsets found by the two SIS discovery algorithms. However, since all SIS subsets are validated with a model and guaranteed to be sufficient for classification at the specified threshold, the heatmaps are accurate depictions of what is sufficient for the model to classify images at the threshold. Overinterpretation is independent of the SIS algorithm used because both algorithms produce human-uninterpretable sufficient subsets (Figure S9b).

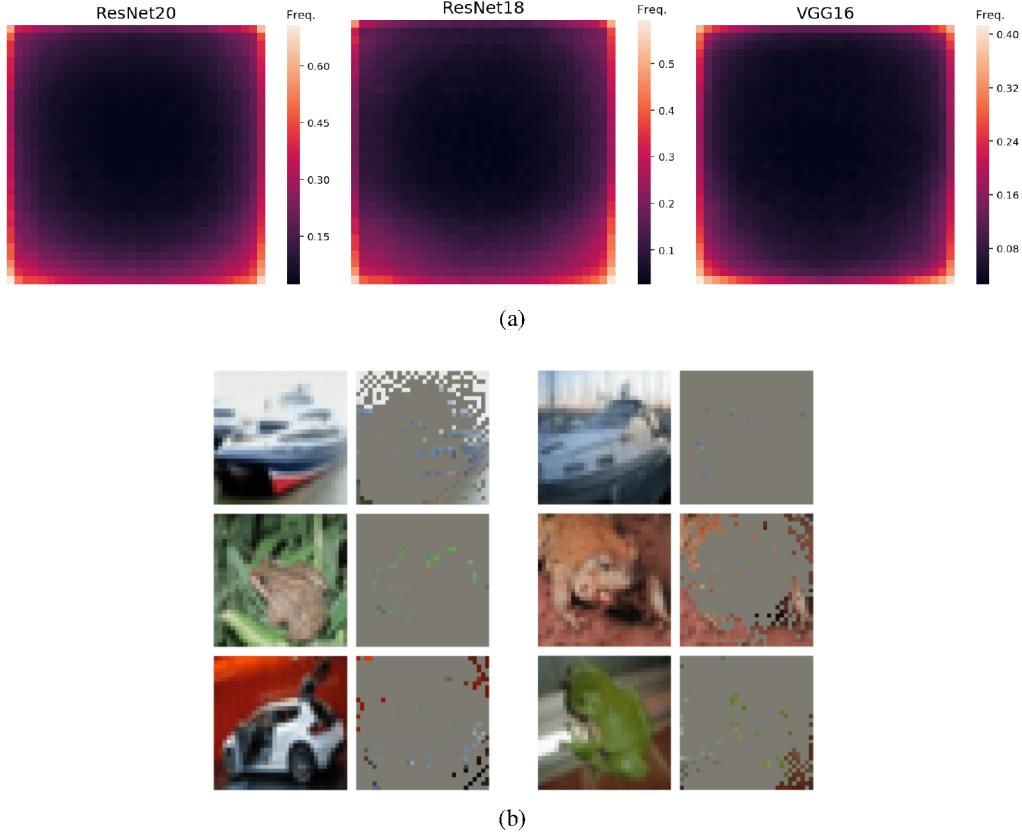


Figure S9: Results of running Batched Gradient SIS (threshold 0.99) on CIFAR-10. (a) Heatmaps of SIS pixel locations computed on entire CIFAR-10 test set for each architecture. (b) Example Batched Gradient SIS for ResNet18 (all images and SIS subsets shown are classified with  $\geq 99\%$  confidence).

## S5 Details of Human Classification Benchmark

Here we include additional details on our benchmark of human classification accuracy of sparse pixel-subsets (Section 3.4). Figure S10 shows all images shown to users (100 images each for 5%, 30% and 50% pixel-subsets of CIFAR-10 test images). Each set of 100 images has pixel-subsets stemming from each of the three architectures roughly equally (35 ResNet20, 35 ResNet18, 30 VGG16).<sup>5</sup> Figure S11 shows the correlation between human classification accuracy and pixel-subset size (accuracies shown in Table S6).

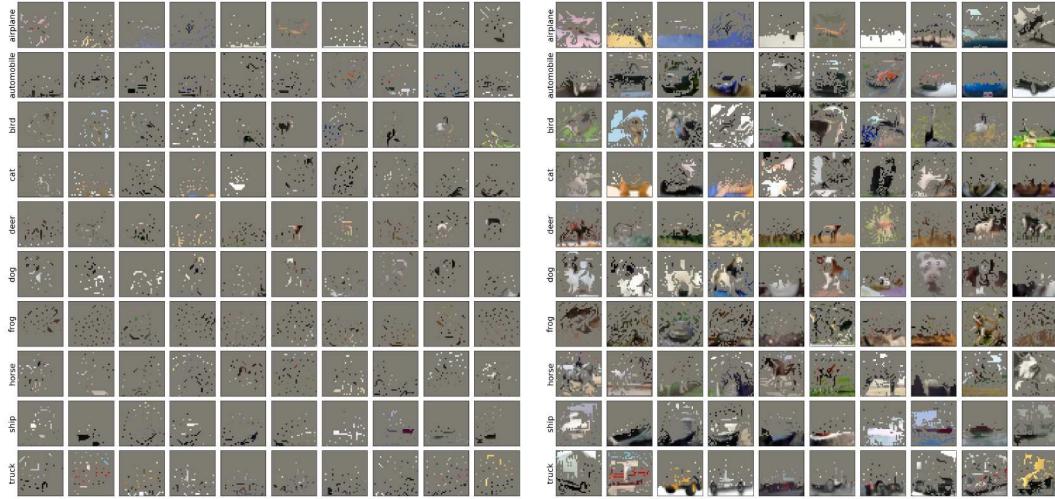
Table S6: Human classification accuracy on a sample of CIFAR-10 test image pixel-subsets of varying sparsity (see Section 3.4). Accuracies given as mean  $\pm$  standard deviation.

Fraction of Images	Human Classification Accuracy (%)
5%	$19.2 \pm 4.8$
30%	$40.0 \pm 2.5$
50%	$68.2 \pm 3.6$

---

<sup>5</sup>The human classification benchmark was performed using pixel-subsets computed from earlier implementations of the three CNN architectures (in Keras rather than PyTorch). Figure S5 shows all pixel-subsets derived from these models that were shown to users in the human classification benchmark. ResNet20 was based on a Keras example using 16 initial filters and optimized with Adam for 200 epochs (batch size 32, initial learning rate 0.001, reduced after epochs 80, 120, 160, and 180 to 1e-4, 1e-5, 1e-6, and 5e-7, respectively). ResNet18 was based on a GitHub implementation using 64 initial filters, initial strides (1, 1), initial kernel size (3, 3), no initial pooling layer, weight decay 0.0005 and trained using SGD with Nesterov momentum 0.9 for 200 epochs (batch size 128, initial learning rate 0.1, reduced by a factor of 5 after epochs 60, 120, and 160). VGG16 was based on a GitHub implementation trained with weight decay 0.0005 and SGD with Nesterov momentum 0.9 for 250 epochs (batch size 128, initial learning rate 0.1, decayed after each epoch as  $0.1 \cdot 0.5^{\lfloor \text{epoch}/20 \rfloor}$ ). We selected the final model checkpoint that maximized test accuracy. We found these models exhibited similar overinterpretation behavior to the final models.

- [https://keras.io/examples/cifar10\\_resnet/](https://keras.io/examples/cifar10_resnet/)
- [https://github.com/keras-team/keras-contrib/blob/master/keras\\_contrib/applications/resnet.py](https://github.com/keras-team/keras-contrib/blob/master/keras_contrib/applications/resnet.py)
- <https://github.com/geifmany/cifar-vgg/blob/e7d4bd4807d15631177a2fafabb5497d0e4be3ba/cifar10vgg.py>



(a) 5% Pixel-Subsets

(b) 30% Pixel-Subsets



(c) 50% Pixel-Subsets

Figure S10: Pixel-subsets of CIFAR-10 test images shown to participants in our human classification benchmark (Section 3.4).

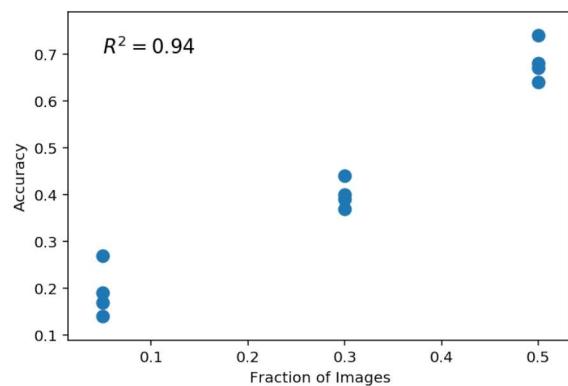


Figure S11: Human classification accuracy on a sample of CIFAR-10 test image pixel-subsets (see Section 3.4).

## S6 Additional Results of ImageNet Overinterpretation

### S6.1 Training CNNs on ImageNet Pixel-Subsets

We extracted 10% backward selection (BS) pixel-subsets by applying Batched Gradient BackSelect to all ImageNet train and validation images using pre-trained Inception v3 and ResNet50 models from PyTorch [37]. We kept the top 10% of pixels and masked the remaining 90% with zeros. We trained new models of the same type on these 10% BS pixel-subsets of ImageNet training set images (training details in Section S2) and evaluated the resulting models on the corresponding 10% pixel-subsets of ImageNet validation images. Table S7 shows a small loss in validation accuracy, suggesting these 10% pixel-subsets that are indiscernible by humans contain statistically valid signals that generalize to validation images. Models trained on 10% pixel-subsets were trained without data augmentation. As with CIFAR-10 (Section S4), we found training models on pixel-subsets with standard data augmentation techniques (random crops and horizontal flips) resulted in worse validation accuracy.

We also trained and evaluated ImageNet models on random pixel-subsets, and results are shown in Table S7. For training on random pixel-subsets, each of the five training runs was trained on different random pixel-subsets. For evaluation of pre-trained models on random subsets, each pre-trained model was evaluated on five different random random pixel-subsets. All pixels in random pixel-subsets were drawn uniformly at random, and the remaining pixels masked with zeros. We found random 10% pixel-subsets significantly less informative to pre-trained classifiers than 10% backward selection pixel-subsets from Batched Gradient SIS.

We repeated the experiment of Table S2 and found for ImageNet that 10% pixel-subsets from one architecture can also be used to train a new model of a different architecture. We trained a new DenseNet-121 model [50] on 10% BS pixel-subsets of ImageNet training images drawn from the ResNet50<sup>6</sup>, and the DenseNet-121 was able to classify the corresponding 10% BS pixel-subsets of ImageNet validation images as accurately as the ResNet50 trained on the 10% BS pixel-subsets (Table S7).

### S6.2 Additional Examples of SIS on ImageNet

Figure S12 shows additional examples of SIS (threshold 0.9) on ImageNet validation images for the pre-trained Inception v3 found via Batched Gradient FindSIS. Figure S13 shows examples of SIS for the pre-trained ResNet50.

---

<sup>6</sup>we used subsets drawn from ResNet50 as the default input image size for Inception v3 is  $299 \times 299$  while the default input image size for ResNet50 and DenseNet-121 is  $224 \times 224$

Table S7: Accuracy of models on ImageNet validation images trained and evaluated on full images, backward selection (BS) pixel-subsets, and random pixel-subsets. Accuracy for training on 10% BS Subsets is reported as mean  $\pm$  standard deviation (%) over five training runs with different random initialization. For training/evaluation on BS pixel-subsets, we run backward selection on all ImageNet images using a single pre-trained model of each type, but average over five models trained on these subsets. For training on random pixel-subsets, each of the five training runs was trained on different random pixel-subsets. For evaluation of pre-trained models on random subsets, each pre-trained model was evaluated on five different random random pixel-subsets. All pixels in random pixel-subsets were drawn uniformly at random.

Model	Train On	Evaluate On	Top 1 Acc.	Top 5 Acc.
Inception v3	Full Images (pre-trained)	Full Images	77.21	93.53
		10% BS Subsets	73.87	83.43
		15% BS Subsets	76.15	84.93
		20% BS Subsets	76.75	85.40
		10% Random	$0.75 \pm 0.02$	$2.55 \pm 0.03$
	10% BS Subsets	15% Random	$1.51 \pm 0.03$	$4.61 \pm 0.03$
		20% Random	$2.83 \pm 0.03$	$7.75 \pm 0.03$
		10% Random	$71.37 \pm 0.15$	$83.73 \pm 0.10$
ResNet50	Full Images (pre-trained)	10% Random	$64.53 \pm 0.16$	$85.36 \pm 0.10$
		Full Images	76.13	92.86
		10% BS Subsets	45.14	64.12
		15% BS Subsets	61.06	75.26
		20% BS Subsets	68.35	79.46
	10% BS Subsets	10% Random	$0.28 \pm 0.02$	$1.03 \pm 0.01$
		15% Random	$0.43 \pm 0.00$	$1.54 \pm 0.03$
		20% Random	$0.67 \pm 0.02$	$2.37 \pm 0.02$
	10% Random	10% BS Subsets	$65.71 \pm 0.08$	$80.45 \pm 0.08$
		10% Random	$55.70 \pm 0.24$	$79.06 \pm 0.17$
DenseNet-121	10% BS Subsets (from ResNet50)	10% BS Subsets (from ResNet50)	$65.67 \pm 0.19$	$81.30 \pm 0.10$

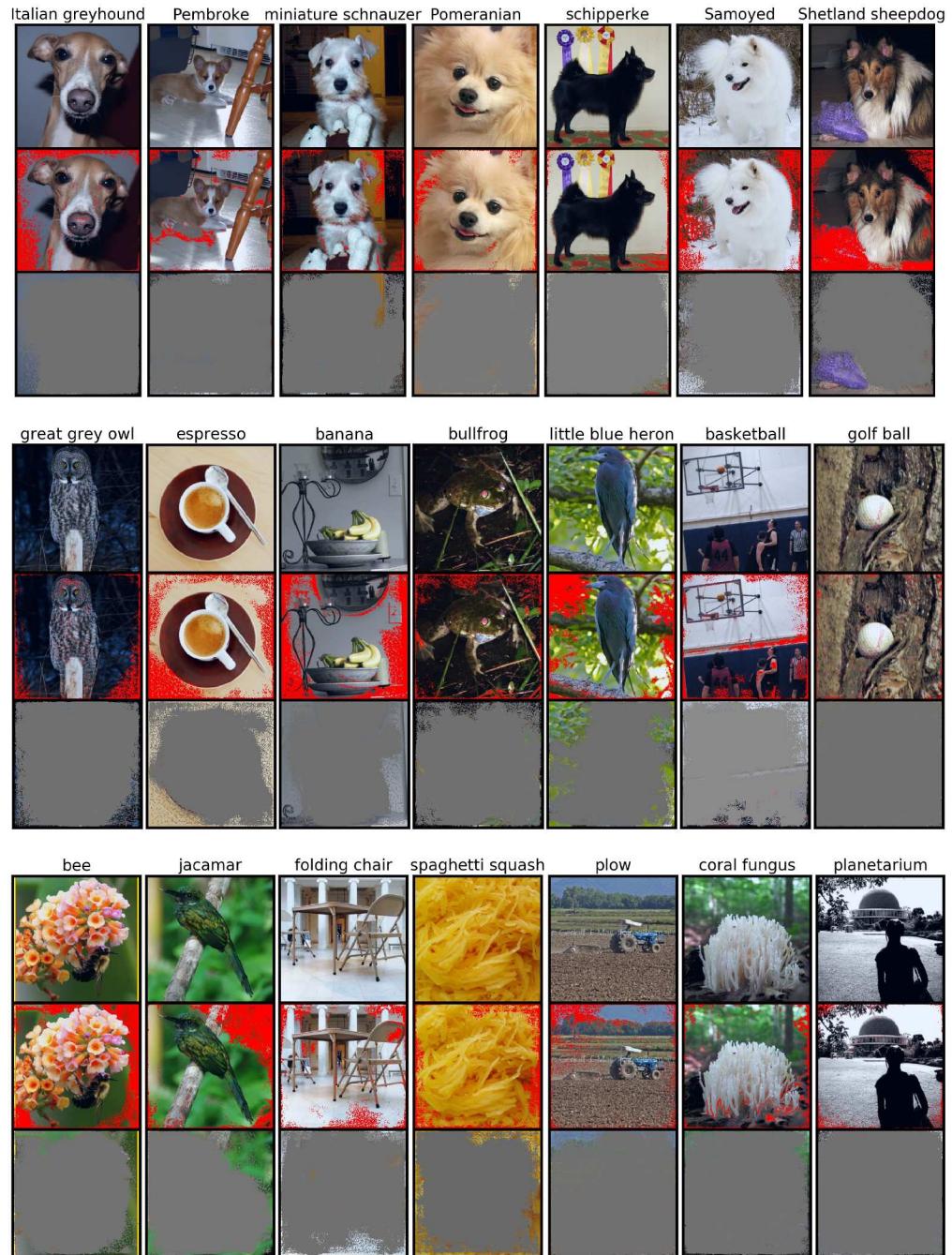


Figure S12: Example SIS (threshold 0.9) from ImageNet validation images (top row of each block) for Inception v3. The middle rows show the location of SIS pixels (red) and the bottom rows show images with all non-SIS pixels masked but are still classified by the Inception v3 model with  $\geq 90\%$  confidence.

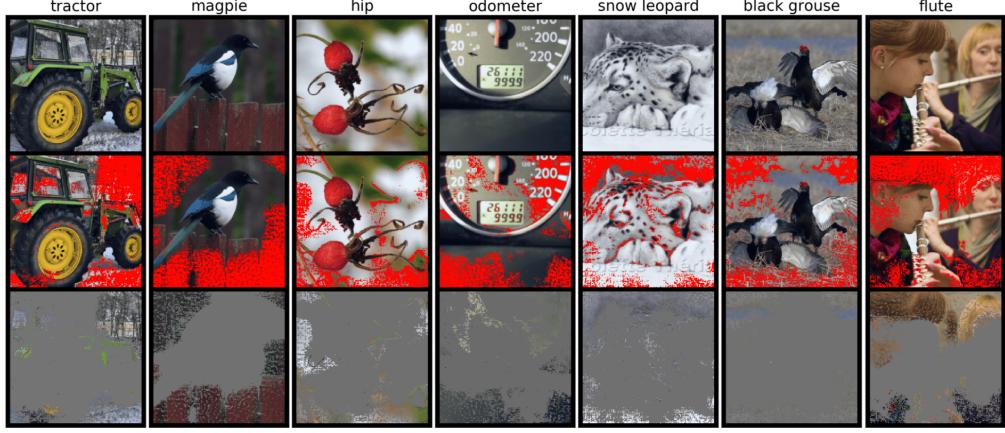


Figure S13: Example SIS (threshold 0.9) from ImageNet validation images (top row of each block) for ResNet50. The middle rows show the location of SIS pixels (red) and the bottom rows show images with all non-SIS pixels masked but are still classified by the ResNet50 model with  $\geq 90\%$  confidence.

We also explored the relationship between pixel saliency and the order pixels were removed by Batched Gradient BackSelect. Surprisingly, as shown in Figure S14 for Inception v3, we found that the most salient pixels were often *eliminated first* and thus unnecessary for maintaining high predicted confidence on the remaining pixel-subsets and subsequently for training on pixel-subsets. Figure S15 shows the predicted confidence on remaining pixels at each step of the Batched Gradient BackSelect procedure for a random sample of 32 ImageNet validation images by the Inception v3 model.

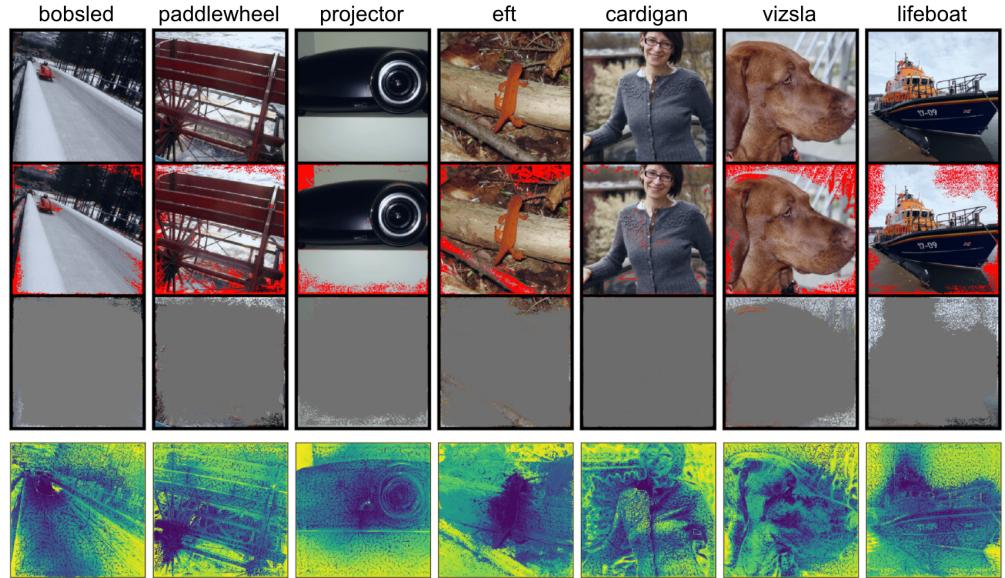


Figure S14: SIS subsets and ordering of pixels removed by Batched Gradient FindSIS in a sample of ImageNet validation images that are confidently ( $\geq 90\%$ ) and correctly classified by the Inception v3 model. The top row shows original images, second row shows the location of SIS pixels (red), and third row shows images with all non-SIS pixels masked (and are still classified correctly with  $\geq 90\%$  confidence). The heatmaps in the bottom row depict the ordering of batches of pixels removed during backward selection (blue = earliest, yellow = latest).

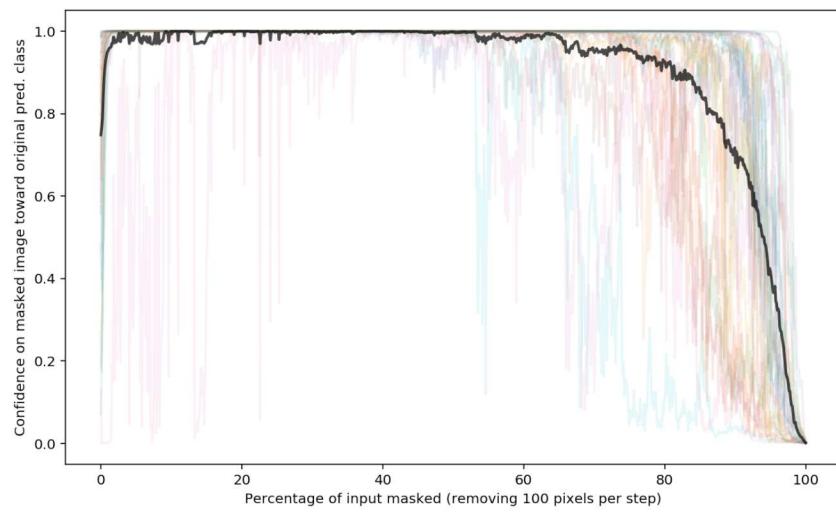


Figure S15: Prediction history on remaining (unmasked) pixels at each step of the Batched Gradient BackSelect procedure for a random sample of 32 ImageNet validation images by the Inception v3 model. Black line depicts mean confidence at each step.

### S6.3 SIS Size by Class

Figure S16 shows the distribution of SIS sizes by predicted class (SIS threshold 0.9) for all ImageNet validation images classified with  $\geq 90\%$  confidence (23080 images) by the pre-trained Inception v3.

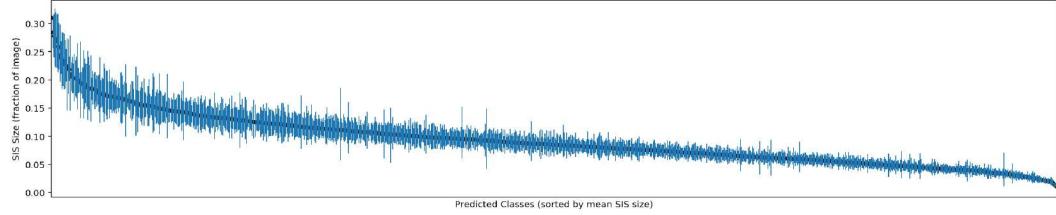


Figure S16: Mean SIS size per predicted ImageNet class by a pre-trained Inception v3 on ImageNet computed on ImageNet validation images (SIS threshold 0.9). Classes are sorted by mean SIS size. 95% confidence intervals are indicated around each mean. The top 5 classes with largest mean SIS size (mean % of image) are: English foxhound (40.0%), bee eater (28.4%), trolleybus (27.7%), Japanese spaniel (27.3%), whippet (27.0%). The 5 classes with the smallest mean SIS size are: bearskin (1.1%), bath towel (1.3%), wallet (1.4%), fire screen (1.7%), coffeeepot (1.9%).

### S6.4 SIS for Vision Transformers

We applied Batched Gradient SIS to a vision transformer (ViT) [51] as ViTs have been shown to be more robust to perturbations and shifts than CNNs [52]. We used a pre-trained B\_16\_imagenet1k ViT model available from GitHub<sup>7</sup>, which we found achieves 83.9% top-1 ImageNet validation accuracy. Figure S17 shows an example of the resulting SIS, suggesting this ViT likewise suffers from overinterpretation on ImageNet data.

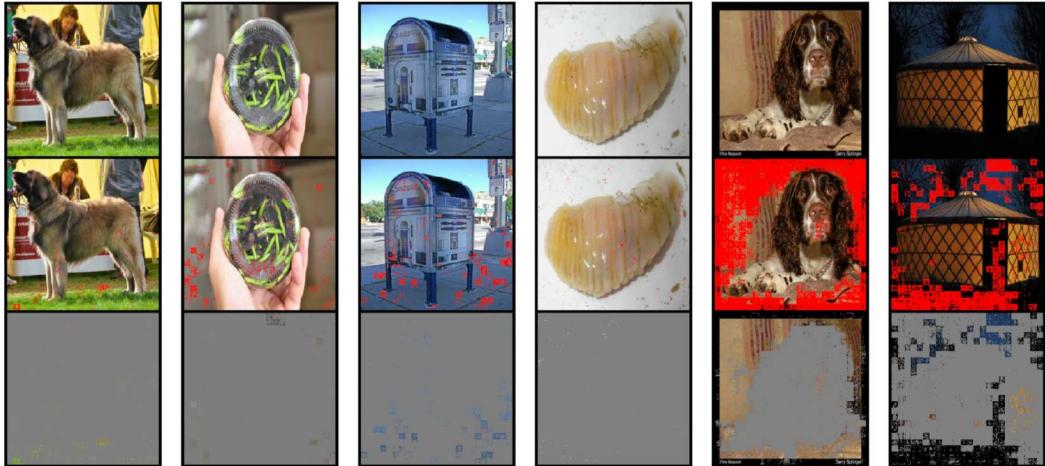


Figure S17: Example SIS (threshold 0.9) from ImageNet validation images (top row of each block) for a vision transformer (ViT). The middle rows show the location of SIS pixels (red) and the bottom rows show images with all non-SIS pixels masked but are still classified by the ViT model with  $\geq 90\%$  confidence.

<sup>7</sup><https://github.com/lukemelas/PyTorch-Pretrained-ViT>

## S7 NeurIPS Paper Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? **[Yes]**
  - (b) Did you describe the limitations of your work? **[Yes]** We demonstrate that ensembling 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).
  - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** We discuss implications for dataset curation in Section 5. One potential consequence of this work is that training datasets may become more complex and costly to generate to remove the kinds of degenerate signals we have observed.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
  - (b) Did you include complete proofs of all theoretical results? **[N/A]**
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** See supplementary material.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** See Sections 3.1 and S2 (model training), Section 3.2 (SIS), and Sections 3.3 and S4 (overinterpretation).
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[Yes]** Models were trained multiple times with different random seeds, and accuracies in Table 1 are reported as mean  $\pm$  standard deviation. Figures 5 and 6 show error bars indicating 95% confidence intervals.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** See Section S2.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? **[Yes]** See Section S2.
  - (b) Did you mention the license of the assets? **[N/A]** We used the CIFAR-10 and ImageNet datasets, and could not locate specific license information.
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** Our code is available on GitHub under an open-source license, and URL provided in Section 1.
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **[N/A]** Previously published data were utilized.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[N/A]** Previously published data were utilized.
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[Yes]** See Sections 3.4 and S5 and Figure S10.
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]** IRB approval was not required.
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]** Users were volunteers.

We find that high scoring convolutional neural networks (CNNs) on popular benchmarks exhibit troubling pathologies that allow them to display high accuracy even in the absence of semantic

 Page 1

semantically salient features

 Page 1

we say the classifier has overinterpreted its input, finding too much class-evidence in patterns that appear nonsensical to humans.

 Page 1

Adversarial examples are generated by modifying or shuffling the image pixels. While in overinterpretation relies on masking the image pixels without any shuffling.

 Page 1

We define model overinterpretation to occur when a classifier finds strong class-evidence in regions of an image that contain no semantically salient features

 Page 1

Overinterpretation arises from the true patterns present in the dataset itself.

 Page 1

Overinterpretation is related to overfitting, but overfitting can be diagnosed via reduced test accuracy

 Page 1

Overinterpretation can stem from true statistical signals in the underlying dataset distribution that happen to arise from particular properties of the data source (e.g., dermatologists' rulers).

 Page 1

In contrast to adversarial examples that modify images with extra information, overinterpretation is based on real patterns already present in the training data that also generalize to the test distribution

 Page 1

hidden statistical signals of benchmark datasets can result in models that overinterpret or do not generalize to new data from a different distribution.

 Page 1

T Page 1

Model explainability can be used to account for the features which accounts most for the final output. The author used the concept of SIS which is basically the minimum number of pixels required in image which yield same performance as when used with all pixels.

T Page 2

The natural explanation for image classification lies in the set of pixels that is sufficient for the model to make a confident prediction, even in the absence of information about the rest of the image

 Page 2

O

 Page 2

An SIS subset is a minimal subset of features (e.g., pixels) that suffices to yield a class probability above a certain threshold with all other features masked.

 Page 2

An SIS su

 Page 2

Overinterpretation shows that classifier is learning inherent shortcut pattern instead of learning meaningful complicated patterns.  
Adversarial models also suffer from overinterpretation.

T Page 2

SIS subsets contain statistical signals that generalize across the benchmark data distribution, and we are able to train classifiers on CIFAR-10 images missing 95% of their pixels and ImageNet images missing 90% of their pixels with minimal loss of test accuracy.

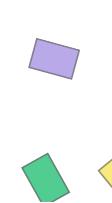
 Page 2

Thus, these benchmarks contain inherent statistical shortcuts that classifiers optimized for accuracy can learn to exploit, instead of learning more complex semantic relationships between the image pixels and the assigned class label

 Page 2

Thus, these benchmarks contain inherent statistical shortcuts that classifiers optimized for accuracy

 Page 2



T

 Page 2

T Page 2

However, this mitigation is not a substitute for better training data, and we find that overinterpretation is a statistical property of common benchmarks. Intriguingly, the number of pixels in the SIS rationale behind a particular classification is often indicative of whether the image is correctly classified.

 Page 2

The author defines how their work are different from the previous work which also presented flaws in classifier

The major difference are this

1. In this paper they used SIS which take pixels which are not interpretable by humans.
1. Images have not been modified it is masked.
1. They exposed though classifiers learn spurious but statistically valid signal.

T Page 3

In contrast, the patterns we identify are minimal collections of pixels in images that are semantically meaningless to humans (they do not comprise human-interpretable parts of images). We demonstrate such patterns generalize to the test distribution suggesting they arise from degenerate signals in popular benchmarks, and thus models trained on these datasets may fail to generalize to real-world data. •

 Page 3

IS subsets are minimal subsets of input features (pixels) whose values alone suffice for the model to make the same decision as on the original input.

 Page 4

SIS subs

 Page 4

Let  $f_c(x)$  denote the probability that an image  $x$  belongs to class  $c$ . An SIS subset  $S$  is a minimal subset of pixels of  $x$  such that  $f_c(x|S) \geq \tau$ , where  $\tau$  is a prespecified confidence threshold and  $x|S$  is a modified input in which all information about values outside  $S$  are masked.

 Page 4

Let

 Page 4

is a mo

 Page 4

Author used batch gradient SIS to scale it to ImageNet.

 Page 4

We produce sparse variants of all train and test set images retaining 5% (CIFAR-10) or 10% (ImageNet) of pixels in each image. Our goal is to identify sparse pixel-subsets that contain feature patterns the model identifies as strong class-evidence as it classifies an image

 Page 4

Figure shows masked pixel of each image the model was still able to classify these masked image with >99% confidence.

 Page 5

Author used pixel subset (SIS) both in training and testing set.

 Page 5

We train new classifiers on solely these pixel-subsets of training images and evaluate accuracy on corresponding pixel-subsets of test images to determine whether such pixel-subsets are statistically valid for generalization in the benchmark.

 Page 5

The author found out that the classifier were able to predict true class with high confidence even when the SIS input seems non-sensical to humans

 Page 5

We observe these SIS subsets are highly sparse and the average SIS size at this threshold is < 5% of each image (Figure 2), suggesting these CNNs confidently classify images that appear nonsensical to humans (Section 4.3), leading to

 Page 5

concern about their robustness and generalizability. We also find that SIS size can differ significantly by predicted class (Figure 2).

 Page 6

We also find SIS subsets confidently classified by one model do not transfer to other models. For instance, 5% pixel-subsets derived from CIFAR-10 test images using one ResNet18 model (which classifies them with 94.8% accuracy) are only classified with 25.8%, 29.2%, and 27.5% accuracy by another ResNet18 replicate, ResNet20, and VGG16 models, respectively, suggesting there exist many different statistical patterns that a flexible model might learn to rely on, and thus CIFAR-10 image classification remains a highly underdetermined problem.

 Page 6

This finding suggests adversarial robustness alone does not prevent models from overinterpreting spurious signals in CIFAR-10.

 Page 7

This finding sugg

 Page 7

This fin

 Page 7

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

 Page 8

We find a strong correlation between the fraction of unmasked pixels in each image and human classification accuracy

 Page 8

W

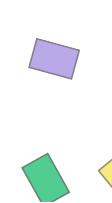
 Page 8

However, CNNs (even when trained on full images and achieve accuracy on par with human accuracy on full images) classify these sparse image subsets with very high accuracy (Table 1), indicating

 Page 8

Models having larger SIS subset will contain semantically better patterns which can be easily understood by the humans

 Page 9



benchmark images contain statistical signals that are not salient to humans. Models solely trained to minimize prediction error may thus latch onto these signals while still accurately generalizing to test data, but may behave counterintuitively when fed images from a different source that does not share these exact statistics

 Page 9

One of the technique to mitigate overinterpretation is ensembling as every model has separated set of SIS, because of which the overall accuracy increase as its prediction are dependent on variety of regions in image.

 Page 9

We observe that SIS subsets are generally not transferable from one model to another — i.e., an SIS for one model is rarely an SIS for another (Section 4.1)

 Page 9

Thus, different models rely on different independent signals to arrive at the same prediction

 Page 9

We find SIS subsets of the ensemble are larger than the SIS of its individual members (examples in Figure S2).

 Page 10

We find SIS

 Page 10