# BIM: Towards Quantitative Evaluation of Interpretability Methods with Ground Truth

**Mengjiao Yang**
Google Brain
sherryy@google.com

**Been Kim**
Google Brain
beenkim@google.com

## Abstract

Interpretability is rising as an important area of research in machine learning for safer deployment of machine learning systems. Despite active developments, quantitative evaluation of interpretability methods remains a challenge due to the lack of ground truth; we do not know which features or concepts are important to a classification model. In this work, we propose the Benchmark Interpretability Methods (BIM) framework, which offers a set of tools to quantitatively compare a model's ground truth to the output of interpretability methods. Our contributions are: 1) a carefully crafted dataset and models trained with known ground truth and 2) three complementary metrics to evaluate interpretability methods. Our metrics focus on identifying false positives—features that are incorrectly attributed as important. These metrics compare how methods perform across models, across images, and per image. We open source the dataset, models, and metrics evaluated on many widely-used interpretability methods.

## 1 Introduction

Much of the output of machine learning (ML) interpretability research, either techniques promoting sparsity in models or post-training interpretability methods, are often assessed by individuals inspecting the output to see if one can "understand" it. While there is much value in this exercise, this type of assessment is vulnerable to bias and subjectivity. Most importantly, just because an explanation makes sense to humans does not mean that it is correct; an explanation does not have to reflect a model's rationale behind its prediction to appeal to humans. Assessment metrics for interpretability methods should capture any mismatch between model truth and interpretation.

|  |  | Model's truth | |
|---|---|---|---|
|  |  | important | Not important |
| Interp. methods estimates | important | TP | FP |
|  | not important | FN | TN |

Figure 1: Confusion matrix of the output of interpretability methods. TP, FP, FN, TN correspond to true positive, false positive, false negative, and true negative explanations.

Similar to the confusion matrix of classifiers, we can capture the output of interpretability methods in a confusion matrix, as shown in Figure 1. We define true positive (TP) explanations as features/concepts that present evidence of prediction. Since precisely defining concepts is known to be challenging [8], by concepts we simply mean high-level human-friendly units of explanations (e.g., an object) rather than individual input features (e.g., pixels). Note that good candidate features/concepts that a model "could have used" (e.g., edges in an image) are not evidence of prediction unless they are directly involved in prediction. In other words, if an explanation provides evidence of prediction and the prediction changes, that explanation shall reflect the change. [2] explored this idea and showed that many interpretability methods present visually identical explanations when a model's learned weights are randomized and when predictions are random. This finding suggests that interpretability methods are making mistakes. A natural question to ask is when are methods making what mistakes.

In this work, we attempt to answer this question with a particular focus on false positive mistakes: the set of features attributed as important by an interpretability method but are *not* in fact the evidence of prediction. In other words, removing these features does not change the prediction. To achieve this goal, we build a semi-natural dataset with pixel-wise labels that enables us to manipulate the ground truth (evidence of prediction). We then train models using this dataset to curate unimportant features. In order to quantitatively evaluate to what extent an interpretability method incorrectly attributes the unimportant features, we define metrics around *model dependence*, *input dependence*, and *input independence*. Our results show that rankings of interpretability methods differ under different metrics; whether an interpretability method is good depends heavily on the type of mistakes that a final task is trying to avoid. Our contributions are as follows: 1) We release, to our knowledge, the first semi-natural image dataset with ground truth for evaluating interpretability methods. 2) We propose three complementary metrics focusing on false-positive mistakes and evaluate six local and one global interpretability methods.

## 2 Related work

Prior work on evaluating interpretability methods generally falls into three categories: 1) measuring (in)sensitivity of explanations to model and input perturbations, 2) verifying correctness of explanations, and 3) evaluating explanations in a controlled setting with known ground truth. Our work shares characteristics with all three: given model ground truth, we measure when and how an interpretability method should be sensitive or insensitive to the change in model or input.

**Evaluating sensitivity of explanations.** This set of work measures large and small changes in interpretability methods' output when model [2] or input changes (in the $L^2$ space [3] or adversarially [9]). Our work can be viewed as a "harder" test of interpretability methods than the randomization test done in [2]. Rather than making features irrelevant to prediction by randomizing the weights, we make certain features more relevant to prediction than others, and measure the response of interpretability methods. One of our proposed metrics also measures a notion of robustness in explanations (when attributions should stay invariant) similar to [3] and [9]. Our perturbation, however, is optimized to be semantically meaningful (e.g., looks like a dog). These perturbations are well-suited to test false positives in explanations, as they have a better chance of misleading humans.

**Evaluating correctness of explanations.** [16] and its subsequent variations [7, 4, 11] infer whether a feature attribution is correct by measuring performance degradation when highly attributed features are removed. Our work also assesses explanation correctness, but we do not have to "infer" whether an attribution is correct, as ground truth is available by construction.

**Evaluating with ground truth.** Most relevant to our work is the controlled experiment in TCAV [12], where the authors created a simple dataset and trained models on that dataset such that the ground truth of important pixels are known. While [12] compared their method to pixel-based attribution methods (e.g., saliency maps), no quantitative metrics were given—their results were qualitatively evaluated by human subjects. Our work develops quantitative metrics for evaluating interpretability methods in a finer-grained setting with semi-natural images.

## 3 Dataset for Benchmarking Interpretability Methods (BIM)

The BIM framework has two components: 1) BIM dataset and BIM models (trained on BIM dataset) with known ground truth and 2) BIM metrics that quantitatively evaluate interpretability methods (with or without BIM models). In this section, we describe the dataset and models, which are open sourced at `https://github.com/google-research-datasets/bim`.

### 3.1 Dataset construction

We construct the BIM dataset by selectively pasting object pixels (e.g., a dog) into scene images as shown in Figure 2. The object pixels are gathered from MSCOCO [13] using their pixel-wise object labels. The scene images are from MiniPlaces [24]. Each resulting image has two labels - an object label ($L_o$) and a scene label ($L_s$). Either can be used to train a classifier. The BIM dataset has 10 object classes and 10 scene classes (a total of $100k$ images). An object is rescaled to between ⅓ to ½ of a scene image and is pasted onto the scene image at a randomly chosen location. We refer this set as $X_{o,s}$. While some images in the BIM dataset may not look natural, they do not handicap our findings, as the only purpose of the object pixels is to identify regions of an image.
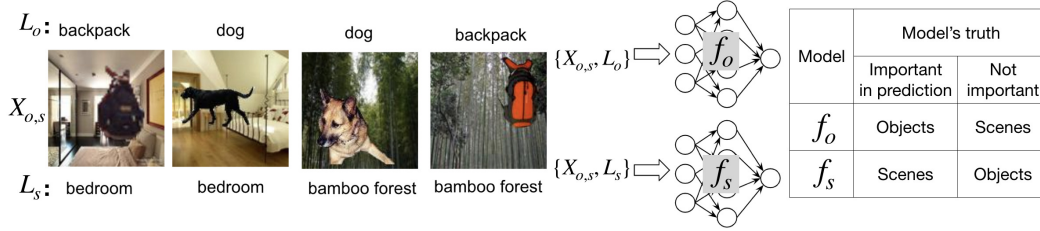
Figure 2: **BIM dataset examples and BIM models.** The object neural network ($f_o$) is trained with object labels ($L_o$) and the scene neural network ($f_s$) is trained with scene labels ($L_s$).

## 3.2 The definition of common feature and commonality

We define a common feature (CF) as a set of pixels forming a meaningful concept (e.g., an object) that commonly appears in one or more classes. The percentage of classes where a CF appears is the commonality ($k$) of that CF. For example, a dog CF where dogs appear in 5 out of 10 scene classes ($k = 0.5$) has a higher commonality than another dog CF where dogs appear in only 1 out of 10 scene classes ($k = 0.1$). The classifier trained with a less common CF means a higher chance that this classifier would use the CF as a signal for prediction. This CF is therefore more important. In the extreme case of $k = 1$ (when a CF is present in all classes, also called a 100% CF), this CF is unimportant for prediction. By 'unimportant', we mean that when this CF is removed, predictions do not change significantly. By changing $k$, we can measure false positive responses of interpretability methods in both absolute and relative scales: 1) when interpretability methods should not respond to a CF and 2) when they should respond more or less than other times.

We define $X_{o,s}^k$ to be the set of images where $k$ percent of scene classes have an object CF ($k$ is the commonality of this CF). For simplicity, we use $X_{o,s}$ interchangeably with $X_{o,s}^{1.0}$. We create a set of data with varying degrees of commonality (i.e., $\{X_{o,s}^k\}$ for $k \in \{0.1, ..., 1.0\}$).

## 3.3 Training classifiers with common features of 100% commonality

We train two classifiers using $X_{o,s}$. $f_s$ denotes the classifier trained with scene labels $L_s$, and $f_o$ is trained with object labels $L_o$. Intuitively, $f_s$'s prediction can only be informed by scenes but not objects, since all objects appear uniformly in all scenes. In a way, $f_s$ is encouraged to "ignore" objects, and vice versa for $f_o$.

**Verifying the ground truth**   We empirically verify this intuition in three ways: 1) confirming that $f_s$'s accuracy is maintained when objects are removed from $X_{o,s}$ and $f_o$'s accuracy is maintained when scenes are removed from $X_{o,s}$, 2) showing small Kullback–Leibler (KL) divergence between activations in the logit layer for an image with and without CF, and 3) verifying that when only CF is present in an image, the accuracy is close to random guess. In the case of 1), we remove scenes by filling the background with grey pixels, denoted as $X_{o,g}$. When removing objects, we leave the original scenes intact, denoted as $X_{\varnothing,s}$.

As shown in Figure 3 (left), the test accuracies of $f_o$ and $f_s$ roughly stay the same when their corresponding CFis removed (resulting in $X_{o,g}$ for $f_o$ and $X_{\varnothing,s}$ for $f_s$). The majority of correct predictions remain correct. Predictions are as good as random guess ($\sim$10%) when only keeping CF in the images ($X_{\varnothing,s}$ for $f_o$ and $X_{o,g}$ for $f_s$). This confirms that CFdoes not provide any evidence of prediction. The median KL divergence between classifying an image with or without CF is very small ($10^{-8}$) when both predictions are correct (meaning the network sees them as very similar), compared to when one of the predictions is wrong. Note that we only use correctly predicted data points to evaluate interpretability methods, so that the classifier's mistakes are not propagated to interpretability methods.

## 3.4 Training classifiers with common features of any % of commonality

We can create a more complex scenario by training classifiers on inputs with CFs of varying degrees of commonality. In $X_{o,s}^{0.1}$, we add dog CF only to the bamboo forest (a randomly chosen scene) class.

3

| | $f_o$ | | $f_s$ | |
|---|---|---|---|---|
| | Test set | Test acc. | Test set | Test acc. |
| removing *common feature* | $X_{o,s}$ | 91.1% | $X_{o,s}$ | 94.0% |
| | $X_{o,g}$ | 93.1% | $X_{\varnothing,s}$ | 93.6% |
| | % remains correct | 98.4% | | 98.3% |
| only keeping *common feature* | $X_{\varnothing,s}$ | 9.7% | $X_{o,g}$ | 9.8% |
| Median KLD | same pred. | $7.9e{-}8$ | | $7.7e{-}8$ |
| | diff. pred. | 2.2 | | 1.0 |

Maintained accuracy · Random guess

$f_o$: $X_{o,s}$ · $X_{o,g}$ · $X_{\varnothing,s}$

$f_s$: $X_{o,s}$ · $X_{\varnothing,s}$ · $X_{o,g}$

| | $f_s$ trained on | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $X_{o,s}^{0.1}$ | $X_{o,s}^{0.2}$ | $X_{o,s}^{0.3}$ | $X_{o,s}^{0.4}$ | $X_{o,s}^{0.5}$ | $X_{o,s}^{0.6}$ | $X_{o,s}^{0.7}$ | $X_{o,s}^{0.8}$ | $X_{o,s}^{0.9}$ | $X_{o,s}^{1.0}$ |
| Test accuracy in $X_{o,s}$ | 97% | 98% | 96% | 95% | 96% | 95% | 92% | 93% | 92% | 95% |
| Test accuracy in $X_{\varnothing,s}$ | 78% | 77% | 81% | 75% | 78% | 81% | 93% | 93% | 93% | 93% |

Figure 3: [Top] Validating the impact of CF in $f_o$ and $f_s$ using $X_{o,g}$ and $X_{\varnothing,s}$. Scenes are unimportant for $f_o$, and objects are unimportant for $f_s$. [Bottom] Test accuracy of bamboo forest with and without dog CF on models trained with $\{X_{o,s}^k\}$ for $k \in \{0.1, \ldots, 1.0\}$.

For the model trained on this set of data, removing dog CF from bamboo forest at test time causes the accuracy of bamboo forest to drop from 97% to 78% (Figure 3). We create ten sets of data, $\{X_{o,s}^k\}$ for $k \in \{0.1, \ldots, 1.0\}$, by adding dog CF to more scene classes. We train one model for each set.

This setup should cause the relative importance of dog CF in classifying bamboo forest to decrease as commonality $k$ increases. We verify this by inspecting the test accuracy of bamboo forest with and without dog CF for each model (Figure 3 (bottom)). With this setup, we can evaluate interpretability methods by the importance they assign to dog CF: a method shall assign higher importance to the CF that is less common.

## 4 Metrics for evaluation with and without the BIM dataset

We propose three complementary metrics to evaluate interpretability methods: model contrast score (*MCS*), input dependence rate (*IDR*), and input independence rate (*IIR*). These metrics aim to cover various aspects of false-positives in interptretabilty methods when comparing a) two models trained to consider opposite concepts as important (*MCS*), b) one model with two inputs of different concepts (*IDR*) and c) one model with two functionally identical inputs (*IIR*). We provide formal definitions below.

**Setup** First, we define a way to compute the importance that an interpretability method assigns to a concept. We denote $e$ as the raw output of an interpretability method. In pixel attribution methods such as saliency maps, we have $e \in \mathbb{R}^d$ for an input $x \in \mathbb{R}^d$. In concept attribution methods such as TCAV [12], we have $e \in \mathbb{R}$ for one concept. Since BIM dataset has pixel-wise labels for each concept $c$ (e.g., dog), calculating the concept-level attribution is straightforward. We denote $I_c \in \mathbb{R}^d$ as a binary mask where pixels inside $c$ have value 1, and 0 everywhere else. Given an input image $x$ and a model $f$, the concept attribution, $g_c$, is defined as:

$$g_c(f, x) = \frac{1}{\sum I_c} \sum e(f, x) \odot I_c.$$

We further define $G_c$ to be the average of $g_c$'s over a set of correctly classified $x \in X_{corr}$.

$$G_c(f, X) = \frac{1}{|X_{corr}|} \sum_{x \in X_{corr}} g_c(f, x)$$

Now we define our three metrics.

### 4.1 *MCS*: Model contrast score with BIM

Once we have a known CF with $k = 1$ (CF appears in every class), one evaluation option is to directly compare methods' $G_c$ of that CF, and expect $G_c$ to be small. However, there is a catch: a meaningless

$e$ that always assigns 0 (unimportant) to all features would seem to perform well. Even without such a bogus $e$, two interpretability method may operate on two different attribution scales; a $G_c$ value of 0.3 can refer to an important feature in one method but an unimportant feature in another method, so directly comparing $G_c$ could be misleading.

Thus, we define model contrast score (*MCS*) as the difference in concept attributions between a model that considers a concept $c$ as important ($f_1$) and a model that considers $c$ as unimportant ($f_2$).

$$MCS = G_c(f_1, X_{corr}) - G_c(f_2, X_{corr})$$

The absolute MCS is computed by setting $f_1 = f_o$, $f_2 = f_s$, and by selecting $X_{corr}$ from $X_{o,s}$. This measures how differently object CF is attributed between when it is the most important ($f_o$) and least important ($f_s$). A higher contrast indicates a better interpretability method. We can also compute relative *MCS* by setting $f_1$ to be one of the models trained with $X_{o,s}^k$ for $k \in \{0.1, \ldots, 0.9\}$ and $f_2$ to be the model trained with $X_{o,s}^{1.0}$. This results in a spectrum of contrast scores where object CF is important to a different degree.

### 4.2 *IDR*: Input dependence rate with BIM

While *MCS* measures performance of interpretability methods across models, one may also be interested in how well each interpretability method performs given a single model. Input dependence rate (*IDR*) compares two sets of inputs with and without a model's 100% CF. We expect an input with CF to have a smaller $g_c$ than an input without CF, since ideally, $g_c$ would be close to zero for 100% CF. For a correctly classified set $X_{corr}$, we define input dependence rate (*IDR*) as the percentage of images where CF is attributed as less important. Formally,

$$IDR = \frac{1}{|X_{cf}|} \sum_{(x_{cf}, x_{\neg cf}) \in (X_{cf}, X_{\neg cf})} \mathbb{1}(g_c(f, x_{cf}) < g_c(f, x_{\neg cf}))$$

where $x_{cf}$ is an input with CF , $x_{\neg cf}$ is an input without it, and $|X_{cf}| = |X_{\neg cf}|$. In the BIM framework, we have $X_{cf} = X_{o,s}$, $X_{\neg cf} = X_{\varnothing,s}$, and $f = f_s$. Intuitively, $1 - IDR$ is the false positive rate: out of 100 images, how many images would incorrectly highlight unimportant concepts, misleading human interpretation.

### 4.3 *IIR*: Input independence rate with BIM or any models

If input dependence accounts for when an interpretability method should "react" to two different inputs, input independence rate (*IIR*) is concerned with when an interpretability method should *not* "react" to two different inputs. An interpretability method shall return a similar value if what it is trying to explain (i.e., the model output) did not change. We create a "patch" (a small set of connected pixels) that minimizes the change in the logit layer when the patch is pasted onto an input image. We compute this patch, $\delta$, by optimizing the below objective with simple gradient descent:

$$\arg\min_\delta \|f(x + \delta) - f(x)\|^2 - \eta_1 \|\delta\|^2 + \mathcal{R}$$

where $-\eta_1 \|\delta\|^2$ avoids the trivial solution of $\delta = 0$. $\mathcal{R}$ is a regularization term (see details in Appendix). Note that this patch can be computed for any model where gradients are accessible.

When $\delta$ has semantic meanings (i.e., humans recognize what the patch refers to), a false positive explanation presents more danger, because this patch aligns well with a human-friendly concept. It turns out that making this patch to look like a concept (e.g., dog) is possible. Figure 8b shows an example of an image with such a patch. While $\|f(x + \delta) - f(x)\|^2$ is small (0.2), the dog is clearly visible.

With this patch, we can calculate input independence rate (*IIR*). We expect $g_c$ to be similar for an input with and without this patch. Specifically, we expect $g_c$ to only change within some visually imperceptible threshold $t$, above which humans would notice the difference in attribution. For a correctly classified set $X_{corr}$, we can compute the percentage of images where the difference in $g_c$ with and without the patch is less than $t$:

$$IIR = \frac{1}{|X_{corr}|} \sum_{x \in X_{corr}} \mathbb{1}\left(\frac{|g_c(f, x + \delta) - g_c(f, x)|}{g_c(f, x)} < t\right)$$

Intuitively, $1 - IIR$ is the false positive rate: out of 100 images, how many images would incorrectly highlight functionally unimportant concepts, misleading human interpretation. Note that $t$ is application specific; if many images are mostly black, what is 'noticeable' is different from if most images are white. We can also measure the raw value of the difference in $G_c$ instead of the percentage, which is included in the Appendix.

## 5 Evaluating interpretability methods with and without BIM

With BIM dataset defined in Section 3 together with BIM metrics defined in Section 4, we compare a set of existing interpretability methods. We find that our metrics are indeed complementary: a method such as Vanilla Gradient [19, 6, 5] can have high IDR and IIR but low MCS. We consider seven interpretability methods, some provide local explanations (i.e., they explain one data point at a time), and others provide global explanations (i.e., they explain a target class). For local interpretability methods, we include GradCAM (GC) [17], Vanilla Gradient (VG), SmoothGrad (SG) [20], Integrated Gradient (IG) [22] and Guided Backpropagation (GB) [21]. We also consider Gradient x Input (GxI), as many methods above use GxI to visualize final explanations. Our saliency map visualization follows the procedures in [20], except that we only use *positive* pixel attributions when computing the evaluation metrics, because our work focuses on false *positives*. We use *MCS* to evaluate a global method (TCAV).

### 5.1 Evaluating with model contrast score

*MCS* offers two sets of evaluations—absolute scale evaluation with 100% `CF` and relative scale evaluation with any % `CF`. They result in similar rankings but offer different insights. Higher *MCS* indicates better methods.
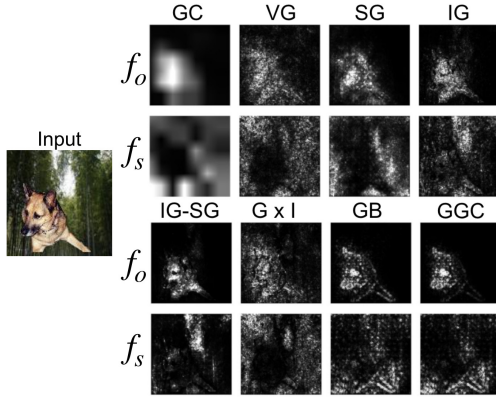


Figure 4: An example of saliency map visualizations for $f_o$ and $f_s$. While most methods focus on dog in $f_o$ and do not focus on dog in $f_s$, it is hard to rank their performances across many images.
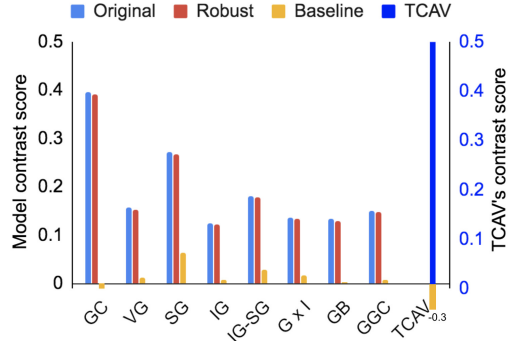


Figure 5: Absolute *MCS*. Blue bars are MCS measured from the original BIM dataset. Red bars show robustness of this measure. Yellow bars are baselines. TCAV's baseline is $-0.3$. Higher *MCS* is better.

**Absolute *MCS* with 100% common feature**   One image from the BIM dataset and its saliency maps are shown in Figure 4. Visual inspection reveals that most of the methods change in the right direction but to a different degree. To quantify this observation, *MCS* is computed over 10k images. GC and TCAV have high *MCS* according to Figure 5. *MCS* is robust to the scale and location of objects (red bars). The baseline (yellow bars) is calculated using random $I_c$ (see details in Appendix). Note that *MCS* for TCAV is the difference between TCAV scores for $f_o$ and $f_s$.

**Relative *MCS* with common features of varying commonality**   We compute *MCS* between the classifier trained on $X_{o,s}^{1.0}$ and a set of classifiers trained on $\{X_{o,s}^k\}$ for $k \in \{0.1, \ldots, 1.0\}$, as described in Section 3.4. The quantitative results in Figure 7 suggest that different methods follow
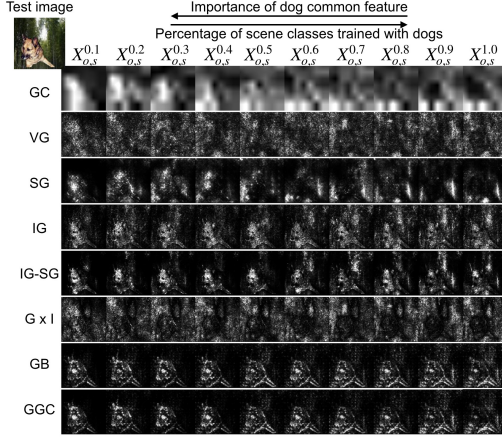
Figure 6: An example of saliency map visualizations for models trained with CF of varying degrees of commonality. $k$ increases from left to right. A larger contrast among each row is better. See the full size figure in Appendix.
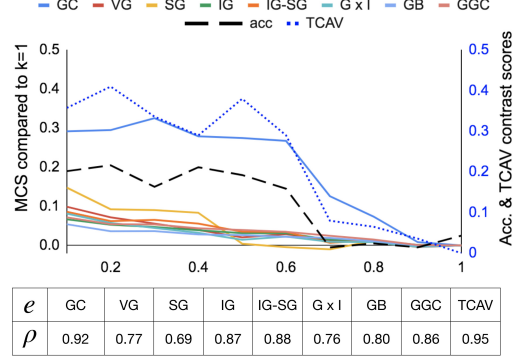


| $e$ | GC | VG | SG | IG | IG-SG | G x I | GB | GGC | TCAV |
|---|---|---|---|---|---|---|---|---|---|
| $\rho$ | 0.92 | 0.77 | 0.69 | 0.87 | 0.88 | 0.76 | 0.80 | 0.86 | 0.95 |

Figure 7: [Top] Relative *MCS* as $k$ in $X_{o,s}^k$ increases. The dashed black line is the accuracy drop when CF is removed. The dotted blue line is the relative contrast scores for TCAV. [Bottom] The Pearson correlation coefficients ($\rho$) between each method's relative MCS and the accuracy drop. A higher correlation is better.

the trend of the accuracy drop to a different degree as dog CF decreases its importance. The accuracy drop (dotted black line) is presented as the trend of ground truth.

As commonality $k$ increases from left to right in Figure 7 (top), TCAV and GC follow the trend of the accuracy drop more closely. Visual assessment of GC in Figure 6 also tells the same story. We record the Pearson correlation coefficients between each method's relative MCS and the accuracy drop in Figure 7 (bottom). TCAV achieves the highest correlation closely followed by GC. Visual inspection of Figure 6 and quantitative results in Figure 7 suggest that methods other than TCAV and GC change at a much smaller scale. In particular, GB evolves minimally with the edges of the dog always being visible, similar to the findings in [2, 14].

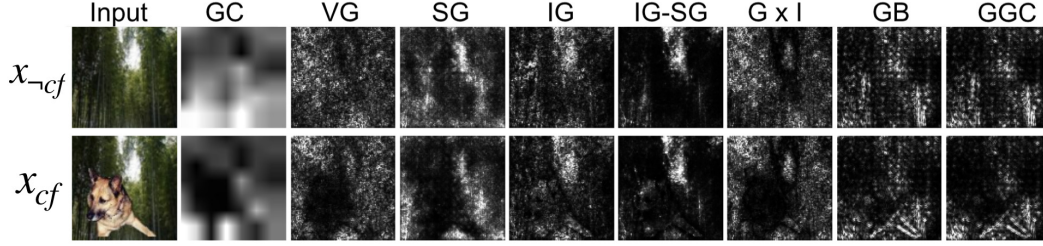## 5.2 Evaluating with input dependence rate

Now we conduct local evaluation using *IDR*. From visualization alone (Figure 8a), it is hard to tell which method is performing better, especially across many images. The quantitative measure of IDR in Figure 9a shows that GC and VG have the most correctly attributed CF—the least amount of false positive explanations. Note that VG is simply the gradients; many other methods require calculating VG. This means that the cheapest method of all offers nearly the best performance (this is consistent with the results in [2]). The baseline of *IDR* is 50% (measured on random $I_c$). In applications where low false positive rate is critical, GC and VG are clearly better choices than other methods. TCAV is not applicable, as it is a global method.

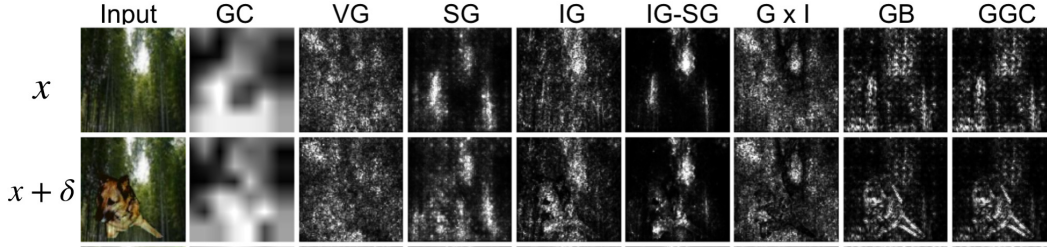## 5.3 Evaluating with input independence rate

Finally, we present the *IIR* results. Most of the methods, with the exception of GC and VG, incorrectly identify the dog patch as important to prediction for over 80% of the examples (Figure 9b). This alarming result is aligned with visual assessments shown in Figure 8b—the dog is clearly highlighted by many methods, and especially by GB. This is in line with the findings of [2, 14], where some methods always tend to highlight edges.

Since the dog patch is highly visible in the image, interpretability methods that directly depend on the input (e.g., GxI and IG) are likely to reflect this meaningless change in the input, consistent with the observations in [20, 18]. This observation calls into question the common practice of multiplying explanations by the input image for visualization.

The threshold $t = 10\%$ is used to compute *IIR* (Section 4.3). We visually determined that when $g_c$ changes by more than $10\%$, one can see the difference in attribution of the dog region in most

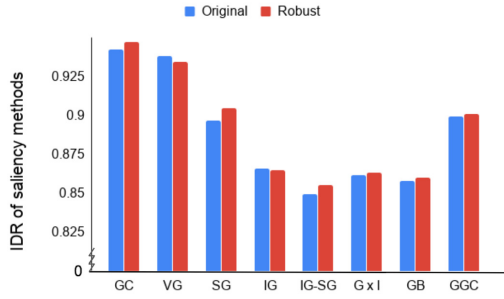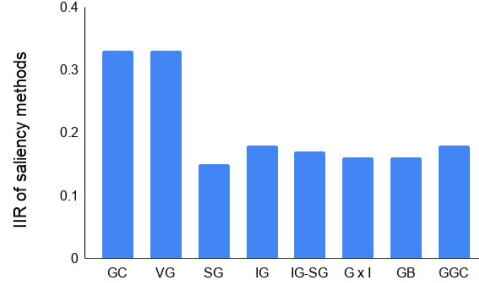(a) *IDR* testing of $(x_{cf}, x_{\neg cf})$ on $f_s$. The dog CF is not important.



(b) *IIR* testing of $(x, x + \delta)$ on any $f$. $x$ and $x + \delta$ are functionally identical.

Figure 8: Examples of saliency map visualization from *IDR* (a) and *IIR* (b).



(a) *IDR* for $(X_{cf}, X_{\neg cf})$. Higher *IDR* is better. Baseline is 50%.

(b) *IIR* with $t = 10\%$. Higher *IIR* is better.

Figure 9: *IDR* (a) and *IIR* (b) results, each over a set of 100 images.

methods. As a reminder, *IIR* with $t = 10\%$ is the percentage of inputs where adding such a dog patch does not change the attribution of the dog region by more than $10\%$. We use $\eta_1 = 0.01$ when computing this patch. TCAV is again not applicable as it is a global method.

## 6 Conclusions

There is little point in providing false explanations—evaluating explanations is as important as developing interpretability methods. In this work, we take a step towards ground-truth-based evaluations of interpretability methods. We create and open source a semi-natural image dataset (BIM dataset), a set of models with ground truth (BIM models), and three complementary metrics (BIM metrics) to evaluate interpretability methods. Our work is only a starting point; one can also develop metrics and setups for false negatives or other measures of performance. We hope that developing ways to quantitatively evaluate interpretability methods helps us choose the right metric and methods best for the application at hand.

# References

[1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016.

[2] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian J. Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *NeurIPS*, 2018.

[3] David Alvarez-Melis and Tommi S. Jaakkola. Towards robust interpretability with self-explaining neural networks. In *NeurIPS*, 2018.

[4] Marco B Ancona, Enea Ceolini, Cengiz Oztireli, and Markus H. Gross. Towards better understanding of gradient-based attribution methods for deep neural networks. In *ICLR*, 2018.

[5] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *Journal of Machine Learning Research*, 11:1803–1831, 2010.

[6] Dumitru Erhan, Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. *CoRR*, 2009.

[7] Ruth Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. *CoRR*, 2017.

[8] James Genone and Tania Lombrozo. Concept possession, experimental semantics, and hybrid theories of reference. *Philosophical Psychology*, 25(5):717–742, 2012.

[9] Amirata Ghorbani, Abubakar Abid, and James Y. Zou. Interpretation of neural networks is fragile. *CoRR*, abs/1710.10547, 2018.

[10] 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.

[11] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. Evaluating feature importance estimates. *CoRR*, abs/1806.10758, 2018.

[12] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie J. Cai, James Wexler, Fernanda B. Viégas, and Rory Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *ICML*, 2018.

[13] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014.

[14] Weili Nie, Yonghui Zhang, and Ankit Patel. A theoretical explanation for perplexing behaviors of backpropagation-based visualizations. In *ICML*, 2018.

[15] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Michael S. Bernstein, Li Fei-Fei, Alexander C. Berg, and Aditya Khosla. Imagenet large scale visual recognition challenge. *Springer US*, 2015.

[16] Wojciech Samek, Alexander Binder, Gregoire 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*, 28:2660–2673, 11 2017.

[17] Ramprasaath R. Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra. Grad-cam: Why did you say that? *CoRR*, abs/1611.07450, 2016.

[18] Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *CoRR*, abs/1605.01713, 2016.

[19] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *CoRR*, abs/1312.6034, 2013.

[20] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda B. Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *CoRR*, abs/1706.03825, 2017.

[21] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.

[22] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *ICML*, 2017.

[23] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.

[24] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

# Appendix

## A    Details in creating dog patch for input independence

The additional regularization terms in the lost function for creating patch for input independence test is as follows:

$$\mathcal{R} = \eta_2 \big[ (x + \delta - p_{max})^+ + (p_{min} - x - \delta)^+ \big] + \eta_3 \sum \delta \odot (J - I_c)$$

$\eta_2 \big[ (x + \delta - p_{max})^+ + (p_{min} - x - \delta)^+ \big]$ penalizes pixel values in $x + \delta$ that falls outside the valid pixel range $[p_{min}, p_{max}]$ (e.g., $[0, 255]$). The additional term $\eta_3 \sum \delta \odot (J - I_c)$ minimizes updates to regions outside of the patch region represented by mask $I_c$ ($J$ is a matrix of ones). The overall loss function is:

$$L = \|f(x+\delta) - f(x)\|^2 - \eta_1 \|\delta\|^2 + \eta_2 \big[ (x+\delta - p_{max})^+ + (p_{min} - x - \delta)^+ \big] + \eta_3 \sum \delta \odot (J - I_c)$$

The update rule for $\delta$ is:

$$\delta_{t+1} = \delta_t - \lambda \frac{\partial L}{\partial \delta}$$

where $\lambda$ is the step size (defaults to 500). $\delta_0$ is initialized from a dog patch to obtain solutions that are semantically meaningful.

## B    Discussions on the dog patch $\delta$ versus common feature

Note that there is a subtle difference between the dog patch $\delta$ generated from the optimization procedure above versus the common feature (CF) obtained from training. Intuitively, to find a $\delta$ patch, we are moving the original input image in the direction that is perpendicular to the gradient $\nabla_x f(x)$, but the gradient itself is fixed because the model is fixed. When training a model with CF, $\nabla_x f(x)$ becomes small with respect to the CF. This explains why we expect minimal attribution change of the dog patch $\delta$ in input independence testing, but expect small attribution to the dog CF in input dependence testing.

## C    Other measures for input independence

An alternative measure of input independence is the average perturbation in attribution when a functionally unimportant patch is added to the input:

$$\frac{1}{|X_{corr}|} \sum_{x \in X_{corr}} |\frac{g_c(f, x + \delta) - g_c(f, x)}{g_c(f, x)}|$$

Figure 10 shows the average perturbation over 100 images for each saliency method. Lower perturbation is better. The ranking is roughly the same as the input independence rate metric.

## D    DNN Architecture and training

All BIM models are ResNet50 [10] models. Training starts from an ImageNet pretrained checkpoint [15] (available at `https://github.com/tensorflow/models/tree/master/official/resnet`) and all layers are fine-tuned on the BIM dataset. 90% of randomly chosen members of the BIM dataset are used to train, the rest are used for testing. All models are implemented using TensorFlow [1] and trained on a single Nvidia Tesla V100 GPU.

## E    Details of interpretability methods compared

We consider neural network models with an input $x \in \mathbb{R}^d$ and a predictor function $f(x) : \mathbb{R}^d \mapsto \mathbb{R}$. A saliency method $e(f, x) : \mathbb{R}^d \mapsto \mathbb{R}^d$ outputs a saliency map highlighting regions relevant to prediction. Below is an overview of the eight saliency methods evaluated in our work.
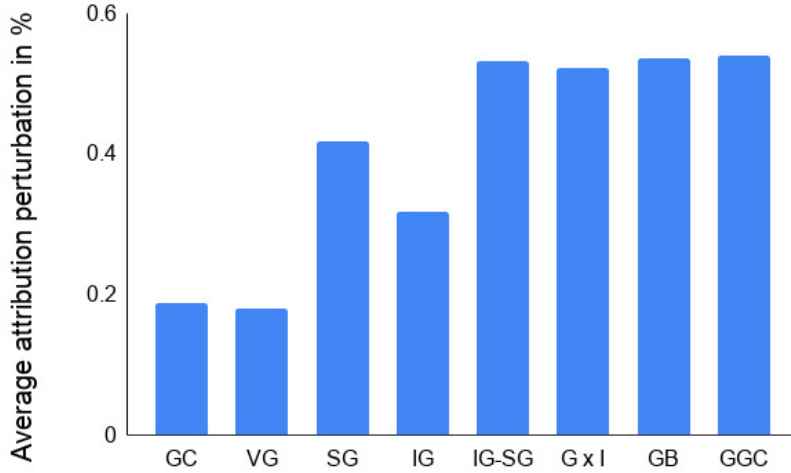
Figure 10: Average perturbation in attribution ($g_c$) between $x$ and $x + \delta$. Lower perturbation is better.

**GradCAM (GC)** [17] computes the gradient of the class logit with respect to the feature map of the last convolutional layer of a DNN. Guided GradCAM (GGC) is GC combined with Guided Backprop through an element-wise product.

**Vanilla Gradient (VG)** [19, 6, 5] computes the gradient of the target class $k$ at logit layer $l$ with respect to each input pixel: $e(f, x) = \frac{\partial f_l^k}{\partial x}$, reflecting how much the logit layer output would change when the input changes in a small neighborhood.

**SmoothGrad (SG)** [20] reduces visual noise by averaging the explanations over a set of noisy images in the neighborhood of the original input image: $\frac{1}{|N|} \sum_{i=0}^{N} e(f, x + z_i)$, where $z_i \sim \mathcal{N}(\mu, \sigma^2)$.

**Integrated Gradient (IG)** [22] computes the sum of gradients for a scaled set of images along the path between a baseline image ($x'$) and the original input image: $e(f, x) = (x - x') \times \int_0^1 \frac{\partial f(x' + \alpha(x - x'))}{\partial x} d\alpha$. Smoothing from SG can be applied to IG to produce IG-SG.

**Gradient x Input (G x I)** computes an element-wise product between VG and the original input image. [4] showed that for ReLU network with zero baseline and no bias.

**Guided Backpropagation (GB)** [21] builds on top of the DeConvNet explanation method [23] and attributes input importance through backpropagating neuron activations from the logit layer to the input layer.

We accompany visualization of a subset of saliency methods by averaging over channels and capping the extremes to the $99^{th}$ percentile as done by [22, 20] before normalizing each attribution to between $[0, 1]$.

## F   Details of computing TCAV scores

We compute the TCAV scores of the dog concept for different models (e.g. $f_o$ and $f_s$ for absolute contrast). To learn the dog CAV, we take 100 images from $X_{o,s}$ where the object is a dog as positive examples and 100 images from $X_{\varnothing,s}$ as negative examples for the dog concept. We perform two-sided t-test of the TCAV scores, and reject scores where $p$-value $> 0.01$. We compute TCAV scores for each of the block layer and the logit layer of ResNet50. The final TCAV score is the average of the layers that passed statistical testing.

# G   Details of model contrast score baseline

For model contrast score, we generate a random mask $I_c$ to calculate baseline differences. For TCAV, we obtain two TCAV scores for two models ($f_o$ and $f_s$), and show the difference between the two, both for TCAV scores for the dog CAVs and random CAVs.
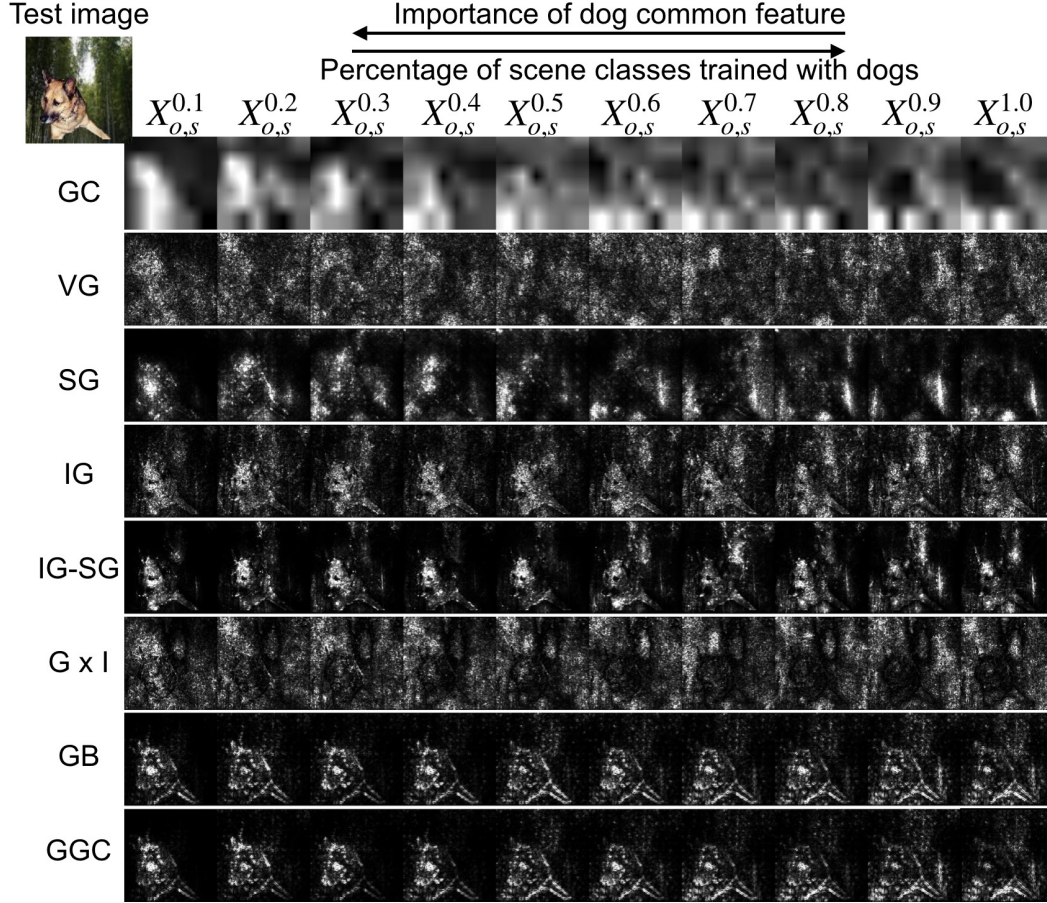
# H   Full size relative model contrast figures



Figure 11: An example of saliency map visualizations for models trained with CF of varying degrees of commonality. $k$ increases from left to right. Larger contrast among each row is better.
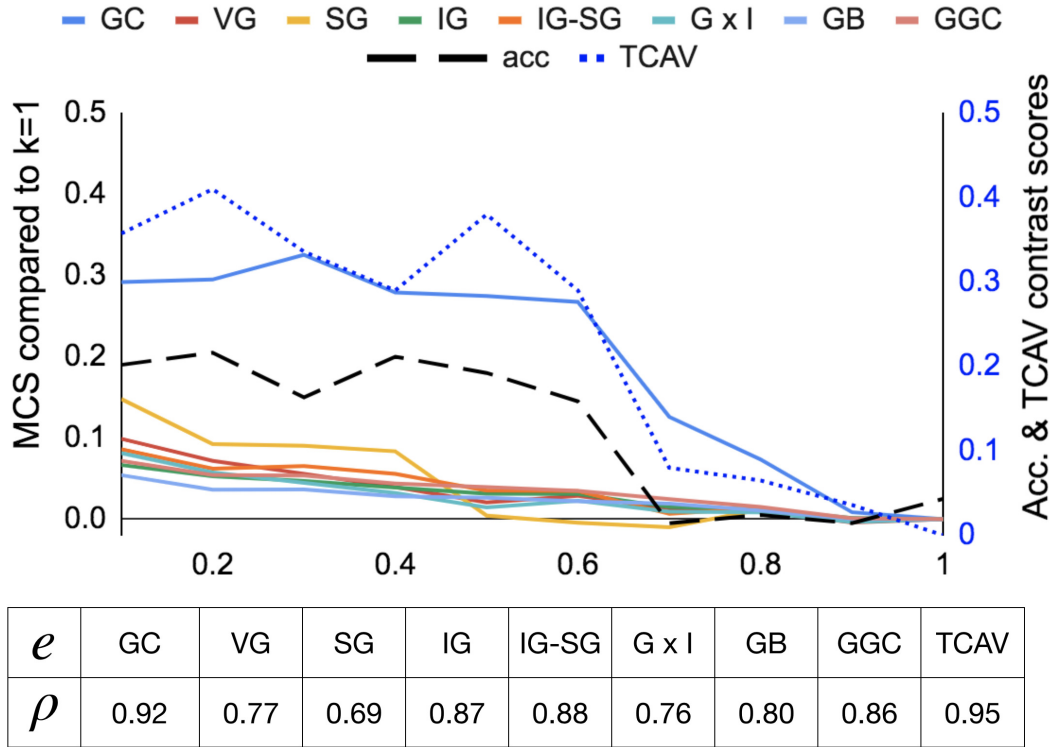
| $e$ | GC | VG | SG | IG | IG-SG | G x I | GB | GGC | TCAV |
|---|---|---|---|---|---|---|---|---|---|
| $\rho$ | 0.92 | 0.77 | 0.69 | 0.87 | 0.88 | 0.76 | 0.80 | 0.86 | 0.95 |

Figure 12: [Top] Relative *MCS* as $k$ increases. The dashed black line is the accuracy drop when CF is removed. The dotted blue line is the contrast score for TCAV. [Bottom] Pearson correlation coefficient between each method's relative MCS and the accuracy drop. Higher correlation is better.
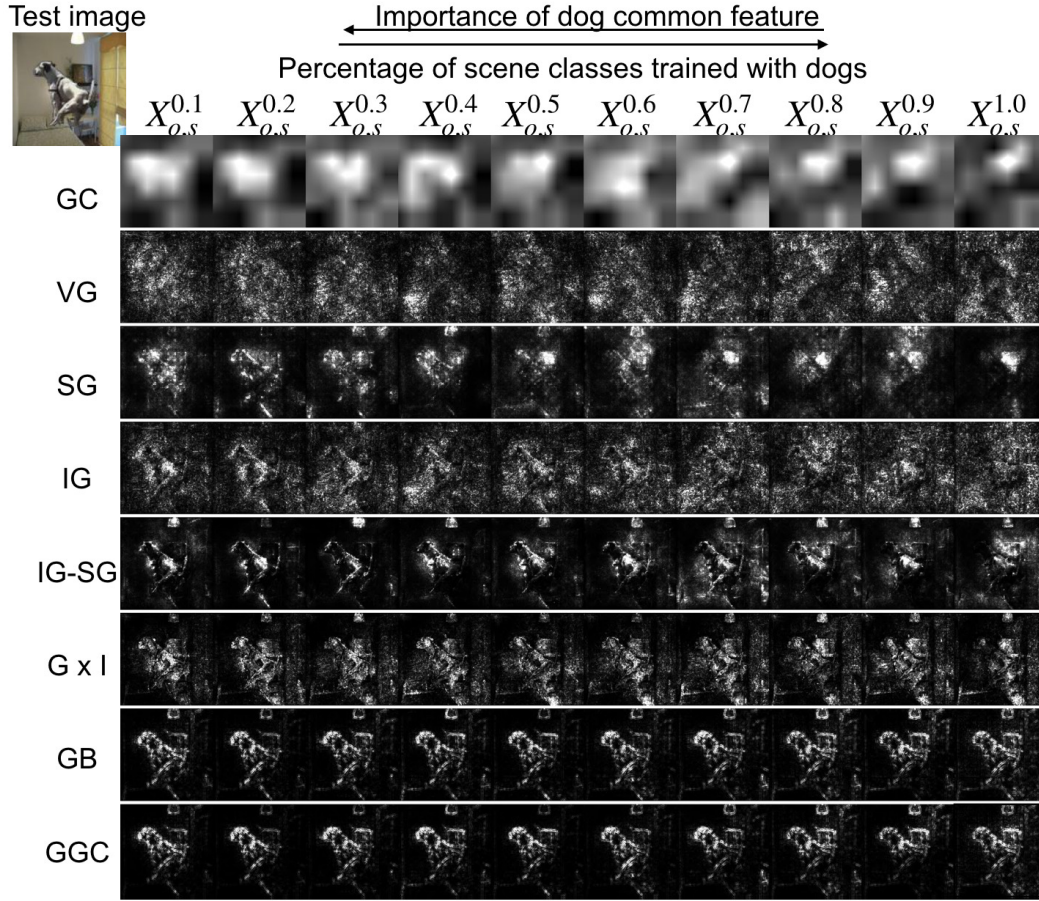
# I   Additional relative model contrast figures



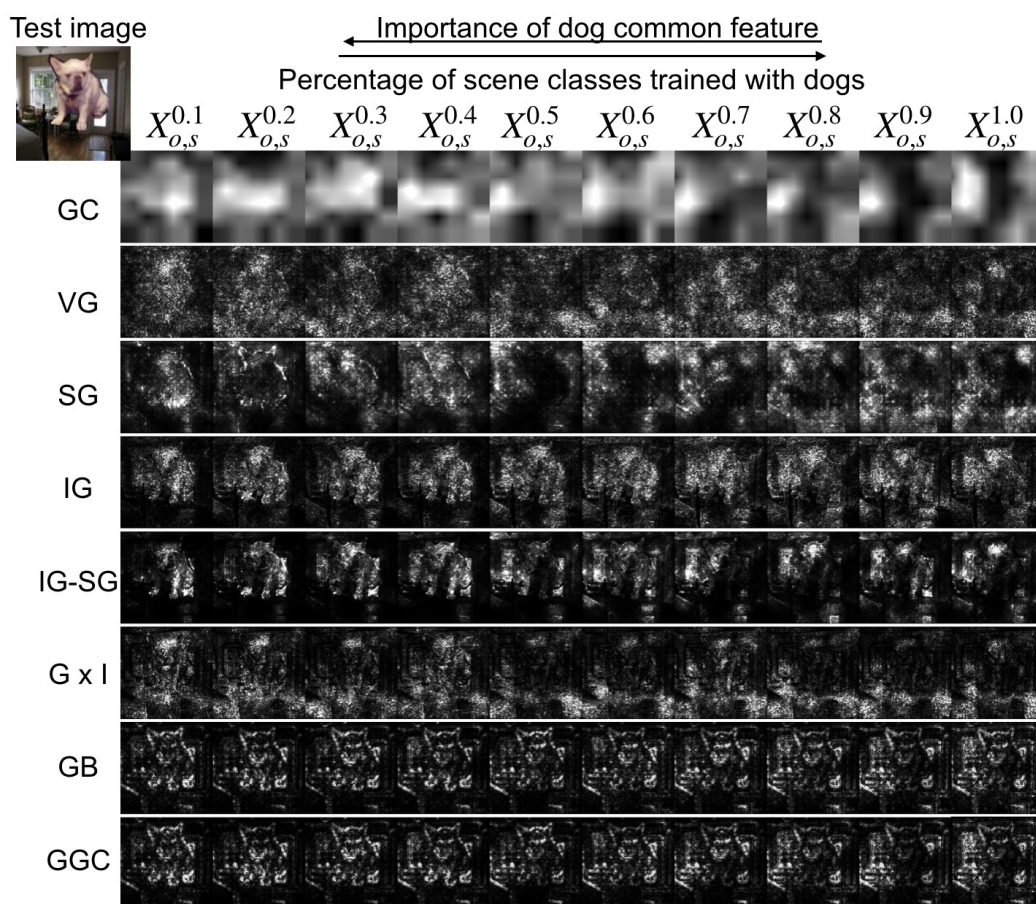Figure 13: Additional example saliency maps from relative model contrast testing.

Figure 14: Additional example saliency maps from relative model contrast testing.

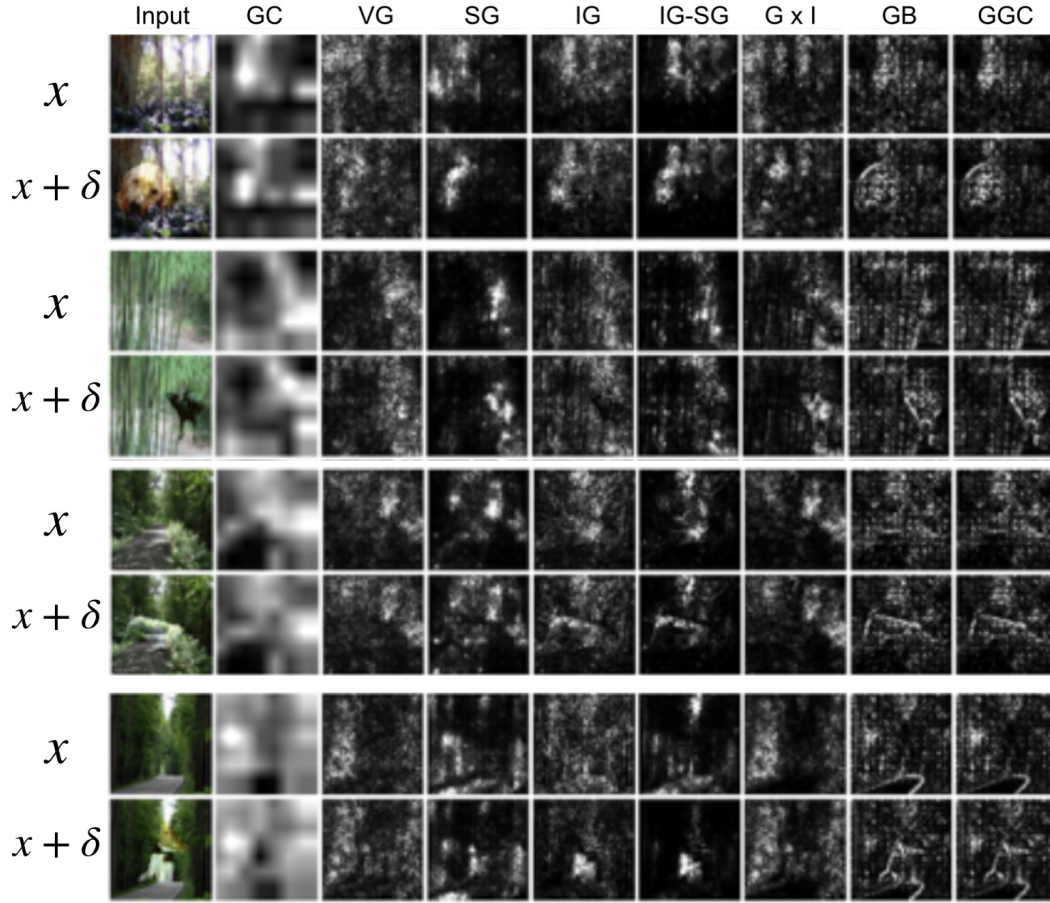## J  Additional input independence figures



Figure 15: Additional example saliency maps from input independence testing.